

# OnePiece: Bringing Context Engineering and Reasoning to Industrial Cascade Ranking System

Sunhao Dai<sup>†1</sup>, Jiakai Tang<sup>‡1</sup>, Jiahua Wu<sup>2</sup>, Kun Wang<sup>2</sup>, Yuxuan Zhu<sup>2</sup>, Bingjun Chen<sup>2</sup>, Bangyang Hong<sup>2</sup>, Yu Zhao<sup>2</sup>, Cong Fu<sup>‡2</sup>, Kangle Wu<sup>2</sup>, Yabo Ni<sup>2</sup>, Anxiang Zeng<sup>2</sup>, Wenjie Wang<sup>3</sup>, Xu Chen<sup>1</sup>, Jun Xu<sup>1</sup> and See-Kiong Ng<sup>4</sup>

<sup>1</sup>Renmin University of China, <sup>2</sup>Shopee, <sup>3</sup>University of Science and Technology of China, <sup>4</sup>National University of Singapore

<sup>†</sup>sunhaodai@ruc.edu.cn, <sup>‡</sup>tangjiakai5704@ruc.edu.cn, <sup>¶</sup>fc731097343@gmail.com

Despite the growing interest in replicating the scaled success of large language models (LLMs) in industrial search and recommender systems, most existing industrial efforts remain limited to transplanting Transformer architectures, which bring only incremental improvements over strong Deep Learning Recommendation Models (DLRMs). From a first principle perspective, the breakthroughs of LLMs stem not only from their architectures but also from two complementary mechanisms: *context engineering*, which enriches raw *input* queries with contextual cues to better elicit model capabilities, and *multi-step reasoning*, which iteratively refines model *outputs* through intermediate reasoning paths. However, these two mechanisms and their potential to unlock substantial improvements remain largely underexplored in industrial ranking systems.

In this paper, we propose OnePiece, a unified framework that seamlessly integrates LLM-style context engineering and reasoning into both retrieval and ranking models of industrial cascaded pipelines. OnePiece is built on a pure Transformer backbone and further introduces three key innovations: (1) structured context engineering, which augments interaction history with preference and scenario signals and unifies them into a structured tokenized input sequence for both retrieval and ranking; (2) block-wise latent reasoning, which equips the model with multi-step refinement of representations and scales reasoning bandwidth via block size; (3) progressive multi-task training, which leverages user feedback chains to effectively supervise reasoning steps during training. Extensive offline experiments verify the effectiveness of each core design and show that OnePiece not only achieves higher sample efficiency but also continues to benefit from larger training spans, surpassing strong baselines with fewer days of logs and maintaining improvement as more data becomes available. OnePiece has been deployed in the main personalized search scenario of Shopee and achieves consistent online gains across different key business metrics, including over +2% GMV/UU and a +2.90% increase in advertising revenue.

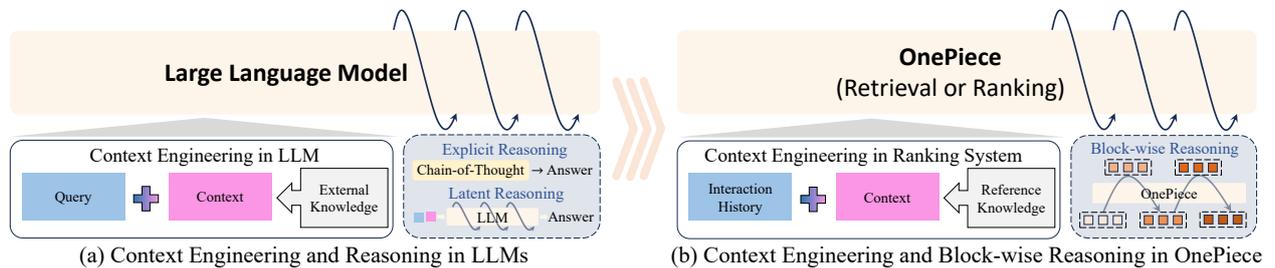


Figure 1 | Bringing LLM mechanisms of context engineering and reasoning to industrial ranking systems. (a) In LLMs, context engineering broadens the query with external knowledge from the input side, while explicit or latent reasoning guides the output process, together enhancing the model’s ability to generate accurate answers. (b) In OnePiece, user interaction history is augmented with additional reference signals, while block-wise reasoning progressively refines representations, jointly improving performance in both retrieval and ranking modes.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Preliminary</b>	<b>4</b>
2.1	Problem Formulation . . . . .	4
2.2	Cascade Ranking Paradigm . . . . .	5
<b>3</b>	<b>The Proposed OnePiece Framework</b>	<b>5</b>
3.1	Overview . . . . .	5
3.2	Context Engineering . . . . .	6
3.3	Backbone Architecture . . . . .	9
3.4	Progressive Multi-Task Training . . . . .	11
3.5	Time Complexity Analysis . . . . .	13
<b>4</b>	<b>Offline Experiments</b>	<b>14</b>
4.1	Experimental Settings . . . . .	14
4.2	Overall Performance . . . . .	15
4.3	Ablation Study . . . . .	16
4.4	Scaling Analysis . . . . .	18
<b>5</b>	<b>Online A/B Testing</b>	<b>19</b>
5.1	Experimental Settings . . . . .	19
5.2	Overall Performance . . . . .	21
5.3	Recall Coverage and Exclusive Contribution . . . . .	21
5.4	Efficiency Analysis . . . . .	22
<b>6</b>	<b>Related Work</b>	<b>23</b>
<b>7</b>	<b>Conclusion and Future Directions</b>	<b>24</b>
	<b>References</b>	<b>25</b>
<b>A</b>	<b>Offline Dataset Construction</b>	<b>30</b>
<b>B</b>	<b>Context Engineering Details</b>	<b>30</b>
<b>C</b>	<b>Attention Visualization Analysis</b>	<b>31</b>

## 1. Introduction

Large Language Models (LLMs), particularly those equipped with reasoning capabilities, have achieved remarkable success across a wide range of complex natural language processing (NLP) tasks (Zhao et al., 2023). Inspired by these advances, researchers have begun exploring how to replicate the key mechanisms behind LLM success in other domains in the hope of realizing similar breakthroughs (Lin et al., 2025; Zhu et al., 2023). In industrial ranking systems (e.g., search and recommender systems), most efforts have focused on transplanting Transformer-based sequential architectures (Wu et al., 2024). While these models bring certain benefits, the improvements over strong Deep Learning Recommendation Models (DLRMs) (Naumov et al., 2019) remain incremental, as core mechanisms such as attention and sequential feature modeling are already deeply integrated into classical designs.

Beyond their foundational Transformer-based architectures, the breakthroughs of LLMs can also be traced to two complementary mechanisms operating at different stages of the modeling process: (1) **Context Engineering** (Dong et al., 2024; Mei et al., 2025), which strengthens the *input side* by enriching raw queries with structured contextual cues to better elicit model capacity; and (2) **Multi-Step Reasoning** (Huang and Chang, 2023; Zhu et al., 2025), which enhances the *output side* by refining predictions iteratively through intermediate latent or explicit reasoning steps. From a first principle view, both context engineering and multi-step reasoning serve the same ultimate goal of eliciting stronger model capacity, one from the *input* side and the other from the *output* side. Together, these two mechanisms substantially expand the generalization ability and capability frontier of LLMs, yet they remain largely underexplored in industrial ranking systems.

However, transplanting these two mechanisms directly into industrial ranking systems is far from straightforward. Unlike LLMs, ranking models cannot readily exploit prompt-style contexts or chain-of-thought supervision, making it unclear how to effectively operationalize context engineering and reasoning in this domain. Specifically, bringing these mechanisms into ranking requires us to tackle two fundamental challenges:

- **How to construct an informative input context?** Most transformer-based industrial models rely primarily on raw user–item interaction sequences, which lack the structural richness of LLM-style prompts (Chen et al., 2025; Deng et al., 2025; Zhai et al., 2024). Moreover, existing feature engineering practices remain predominantly tailored to DLRM-style architectures, leaving open the question of how to enrich context to better endow ranking models with reasoning capabilities.
- **How to optimize multi-step reasoning?** In LLMs, large-scale chain-of-thought annotations provide explicit guidance for training and optimizing reasoning processes (Hao et al., 2024; Shen et al., 2025). In contrast, industrial ranking systems lack such supervision, and even domain experts cannot feasibly articulate the latent decision paths underlying user behaviors in natural language (Tang et al., 2025). This makes it difficult to directly supervise reasoning trajectories.

To bridge this gap, we propose OnePiece, a unified framework that seamlessly integrates LLM-style context engineering and reasoning into both the retrieval and ranking stages of industrial cascaded pipelines. **First**, OnePiece unifies input representation across retrieval and ranking through *structured context engineering*. Specifically, it enriches user interaction history with preference anchors (auxiliary item sequences constructed from domain knowledge provide reference signals, e.g., top-clicked items under the current query) and situational descriptors (e.g., user profiles and query contextual features). For ranking, OnePiece further augments the context with a compact candidate item set, where items are jointly visible to each other, enabling the model to capture cross-candidate interactions for more accurate score prediction. **Second**, OnePiece introduces *block-wise latent reasoning*, where hidden states are progressively refined across multiple reasoning blocks, each building upon the last, thereby expanding the reasoning bandwidth and yielding more expressive representations. **Finally**, to effectively supervise this multi-step process, OnePiece adopts a *progressive multi-task training*

strategy that leverages naturally available user feedback chains (e.g., click, add-to-cart, order) as staged supervision signals. This design provides structured guidance for intermediate reasoning steps and enables the model to gradually align with deeper levels of user preference, from surface engagement to final conversion.

We evaluate the effectiveness of OnePiece through extensive offline experiments and online A/B tests. A series of ablation studies demonstrates the effectiveness of each key module: context engineering improves Recall@100 by +0.110 and click-AUC by +0.109, while block-wise reasoning adds another +0.047 Recall@100 and +0.030 click-AUC. Moreover, OnePiece achieves stronger data efficiency than DLRM and continues to improve as training spans grow longer, while DLRM quickly plateaus. Beyond offline evaluation, OnePiece has been fully deployed in Shopee’s main personalized search scenario, serving billions of users. Online A/B testing confirms that OnePiece delivers consistent business gains: retrieval raises GMV/UU by +1.08%, while ranking yields a +1.12% GMV/UU gain together with a +2.90% increase in advertising revenue. Notably, in retrieval A/B testing, OnePiece covers nearly 70% of impressions recalled by other strategies while delivering 2× higher exclusive contribution than DLRM, demonstrating its ability to both subsume existing recall routes and provide substantial novel impressions and clicks.

Our main contributions can be summarized as follows:

- To the best of our knowledge, this is the first work that explores and deploys context engineering and multi-step reasoning in industrial-scale ranking systems, achieving significant improvements over strong DLRM baselines in both retrieval and ranking tasks.
- We propose OnePiece, a unified framework that introduces both structured context engineering and block-wise latent reasoning into retrieval and ranking stages of cascaded pipelines, together with a progressive multi-task training strategy designed to effectively optimize multi-step reasoning.
- We conduct extensive offline and online evaluations, including large-scale A/B testing in Shopee’s main personalized search scenario. These experiments validate the effectiveness of each design choice, showcase favorable scaling and efficiency properties, and confirm the practicality of deploying OnePiece in real-world industrial environments.

## 2. Preliminary

In this section, we illustrate the problem formulation and typical retrieval-ranking cascade paradigm in industrial ranking systems. Without loss of generality, we focus on the personalized search setting where query-related information is present. In recommendation settings without explicit queries, the query-related component can be omitted.

### 2.1. Problem Formulation

A typical ranking system (e.g., a search or recommender system) aims to return an ordered list of items that the user is most likely to interact with next. Formally, let  $\mathcal{U}$ ,  $\mathcal{V}$ , and  $\mathcal{Q}$  denote the sets of users, items, and queries (in search scenarios), respectively. For each user  $u \in \mathcal{U}$ , we define  $\mathbf{u}$  as the complete user feature representation, encompassing both the user ID and user-side attributes (e.g., age, gender). Similarly, for each item  $v \in \mathcal{V}$ , we denote by  $\mathbf{v}$  the complete item feature representation, comprising the item ID and item-side attributes (e.g., category, shop ID, price). In personalized search, each user query  $q \in \mathcal{Q}$  has a feature representation  $\mathbf{q}$  encoding the query ID and query-related information (e.g., textual embeddings). A user’s historical behavior is modeled as a chronological sequence of interacted items:  $S^u = (v_1^u, \dots, v_t^u, \dots, v_{n_u}^u)$ , where  $v_t^u$  denotes the  $t$ -th item the user  $u$  interacted with and  $n_u$  is the number of interactions for user  $u$ . Given a user’s interaction history up to time  $t$ , denoted as  $S_{1:t}^u = (v_1^u, \dots, v_t^u)$ , along with the current context features  $(\mathbf{u}, \mathbf{v}, \mathbf{q})$ , the

personalized ranking system computes a relevance score for each candidate item  $v \in \mathcal{V}$  through a learned scoring function:

$$\mathcal{F}_\theta(S_{1:t}^u, \mathbf{u}, \mathbf{v}, \mathbf{q}) \mapsto \mathbb{R}, \quad (1)$$

where  $\theta$  represents the model parameters. The system subsequently produces a ranked list  $\pi$  by sorting all candidate items in descending order of their scores. In pure recommendation scenarios without an explicit query, the user query  $\mathbf{q}$  is excluded from the scoring function.

## 2.2. Cascade Ranking Paradigm

Cascade ranking has become the dominant paradigm in industrial large-scale ranking systems due to its ability to balance computational efficiency and ranking quality. This paradigm organizes the decision process into multiple stages (e.g., Retrieval, Pre-ranking, Ranking) in a funnel-like manner: early stages employ lightweight models to rapidly filter out the majority of irrelevant items from massive candidate pools, while later stages apply increasingly sophisticated but computationally expensive models to progressively refined and smaller candidate sets, ultimately producing a high-quality ranked list.

For clarity, we consider a typical two-stage cascade ranking system consisting of:

- **Retrieval Stage.** The goal is to efficiently select a small set of promising candidates  $\mathcal{V}'$  from the full corpus  $\mathcal{V}$ . Industrial systems commonly adopt a dual-tower architecture: a *user tower* transforms user context  $(S_{1:t}^u, \mathbf{u}, \mathbf{q})$  into a dense representation  $\mathbf{h}_u$ , while an *item tower* independently transforms each item  $\mathbf{v}$  into  $\mathbf{h}_v$ . Retrieval is then performed via approximate nearest neighbor (ANN) search based on the similarity between user and item representations, typically computed as inner product or cosine similarity.
- **Ranking Stage.** The goal is to produce a finely-ordered list  $\pi$  from retrieved candidates  $\mathcal{V}'$  by estimating precise preference scores. This stage typically adopts a single-tower architecture, where each candidate  $v \in \mathcal{V}'$  is jointly encoded with the full context  $(S_{1:t}^u, \mathbf{u}, \mathbf{q})$  into a feature sequence. The sequence is processed by a Transformer-based backbone to model rich feature interactions, followed by an MLP head to estimate final scores.

In summary, the retrieval and ranking stages differ in both *objectives* and *inputs*: retrieval focuses on coarse-grained recall with disjoint user/item encoders and the entire item corpus as the search space, while ranking performs fine-grained discrimination with joint encoding of the candidate items.

## 3. The Proposed OnePiece Framework

In this section, we present **OnePiece**, a unified framework that brings LLM-style *context engineering* and *multi-step reasoning* into industrial cascaded ranking systems. We begin by introducing the overall design principles of OnePiece, then describe the structured context construction shared across the retrieval and ranking stages. We further detail the block-wise latent reasoning mechanism and the progressive multi-task optimization method.

### 3.1. Overview

**OnePiece** is a unified framework for industrial cascaded ranking systems that combines *structured context engineering*, *block-wise latent reasoning*, and a *progressive multi-task training strategy*:

- **Context Engineering for Input Token Sequence.** A flexible LLM-style input construction that encodes heterogeneous signals into a unified token sequence, including (1) user interaction history,

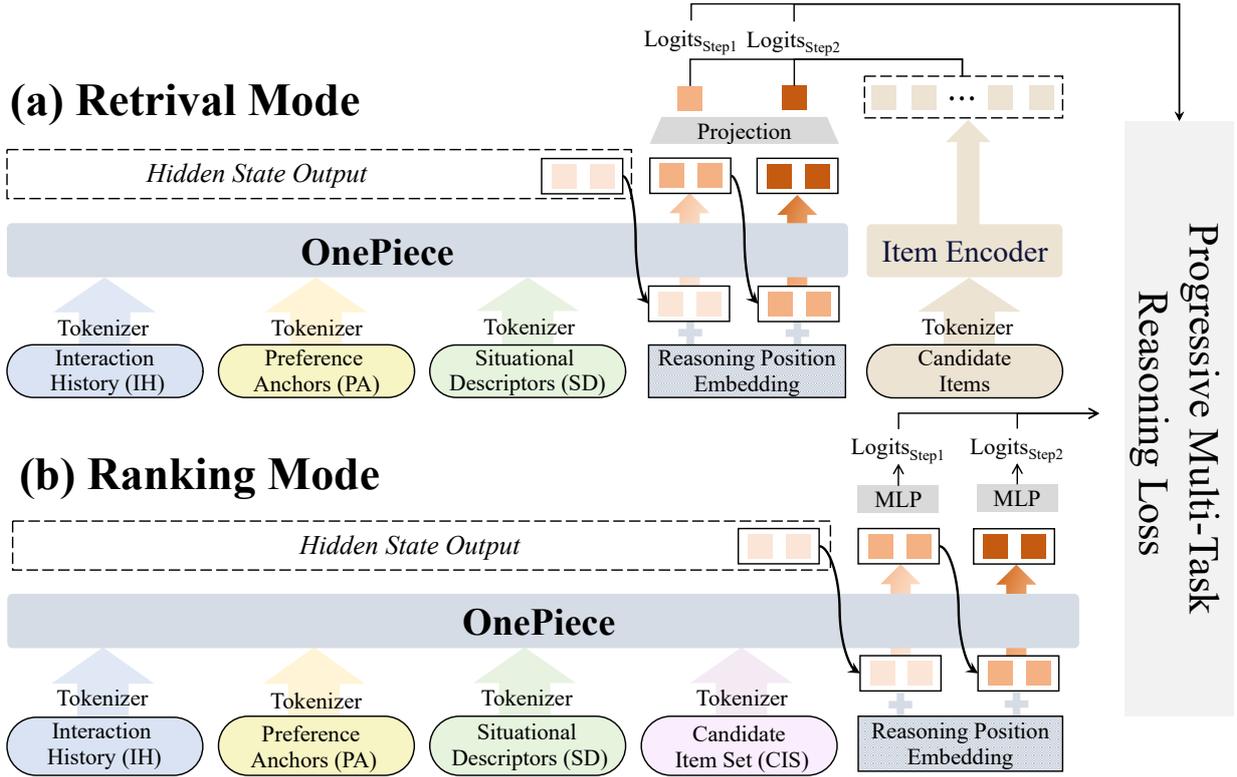


Figure 2 | Overall architecture of the proposed **OnePiece** framework. Retrieval Mode (a) and Ranking Mode (b) both employ structured context engineering to construct unified input tokens, utilize block-wise latent reasoning to iteratively enhance representations across multiple reasoning steps, and are optimized through progressive multi-task training strategy.

(2) preference anchors constructed according to expert experience, (3) situational descriptors capturing scenario-specific context, and (4) candidate item descriptors (used only in ranking).

- **Transformer-based Backbone with Block-wise Reasoning.** A simple and pure Transformer-based backbone augmented with latent reasoning blocks that iteratively refine intermediate representations, enabling step-by-step modeling of user preferences.
- **Progressive Multi-Task Training Strategy.** A staged optimization strategy that supervises different reasoning blocks with multi-level user feedback: earlier blocks are aligned with abundant weak signals (e.g., clicks), while later blocks are guided by stronger but sparser signals (e.g., purchases).

OnePiece unifies retrieval and ranking with a pure Transformer backbone, maintaining the efficiency required in industrial deployment while enhancing context-awareness and reasoning depth across cascade stages. In the following, we will illustrate the details of each component.

### 3.2. Context Engineering

We transform all inputs into a unified *token sequence* that can be directly processed by a transformer-based backbone. To comprehensively capture the user intent and situational context information, we design four complementary token types:

- **Interaction History (IH)** encodes the user’s historical item interactions in chronological order, capturing temporal patterns and evolving preferences.
- **Preference Anchors (PA)** incorporate auxiliary item sequences constructed based on expert knowledge, providing informative cues to capture context-specific user intents beyond the interaction

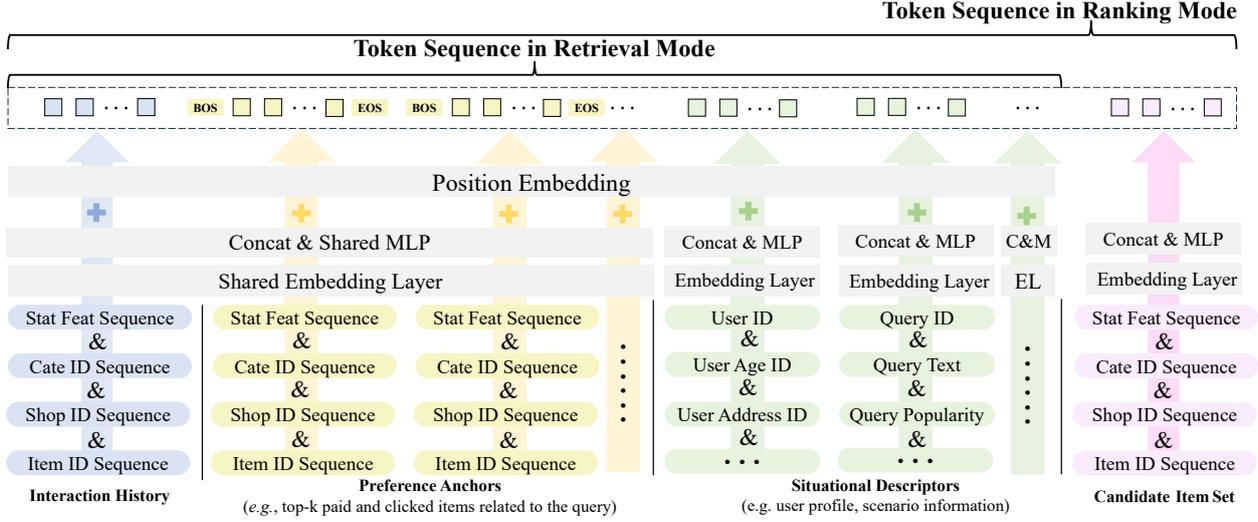


Figure 3 | Context engineering and tokenizer design for input token sequences in **OnePiece**. Both retrieval and ranking share the same construction of interaction history (IH), preference anchors (PA), and situational descriptors (SD). The key difference is that ranking additionally incorporates candidate item set (CIS) tokens, enabling joint scoring within the single-tower architecture.

history.

- **Situational Descriptors (SD)** represent static user features and query-specific information, providing essential context for the current ranking task.
- **Candidate Item Set (CIS, ranking mode only)** contains complete feature representations of candidate items to enable fine-grained comparison and scoring.

This unified sequence representation enables the model to jointly reason over long-term behavioral patterns, domain reference knowledge, and contextual constraints within a unified framework.

Following the preliminaries, we denote a user as  $u$ , a query (when present) as  $q$ , and an item as  $v$ . Their corresponding feature representations are denoted as  $\mathbf{u}$ ,  $\mathbf{q}$ , and  $\mathbf{v}$ , respectively, where each representation contains the identifier and all associated content information. We employ entity-specific embedding functions  $\phi_{\text{user}}(\cdot)$ ,  $\phi_{\text{query}}(\cdot)$ , and  $\phi_{\text{item}}(\cdot)$  to map each entity’s information (including both categorical and continuous features) into a concatenated embedding vector. These embedding functions process different feature types and dimensions according to each entity’s characteristics. To unify these embeddings into the  $d$ -dimensional hidden space of the backbone model, we use lightweight projection heads (e.g., MLPs). Specifically, we define  $\text{Proj}_{\text{user}}$ ,  $\text{Proj}_{\text{query}}$ , and  $\text{Proj}_{\text{cand}}$ , where each projection maps from its corresponding input dimension to  $\mathbb{R}^d$ . The IH and PA components share  $\text{Proj}_{\text{shared}}$ , while user, query, and candidate tokens use their respective projection layers. Next, we detail each component of the input token sequence.

**Interaction History (IH).** Given the user’s chronological interaction sequence  $S^u = (\mathbf{v}_1^u, \dots, \mathbf{v}_{n_u}^u)$ , we embed each item descriptor using a shared projection layer:

$$\mathbf{z}_t^{\text{IH}} = \text{Proj}_{\text{shared}}(\phi_{\text{item}}(\mathbf{v}_t^u)) \in \mathbb{R}^d. \quad (2)$$

We then incorporate temporal information by adding positional embeddings based on the interaction order:

$$\mathbf{h}_t^{\text{IH}} = \mathbf{z}_t^{\text{IH}} + \mathbf{p}_t^{\text{IH}}, \quad 1 \leq t \leq n_u, \quad (3)$$

where  $\mathbf{p}_t^{\text{IH}} \in \mathbb{R}^d$  denotes the learnable positional embedding for the  $t$ -th interaction in the sequence, and  $\mathbf{h}_t^{\text{IH}}$  represents the final token embedding that integrates both content and temporal information.

**Preference Anchors (PA).** Preference anchors are constructed based on domain knowledge to guide the model’s reasoning process. The anchors serve as high-quality reference points that inject inductive biases into the sequence, thereby constraining the representation space and steering the model toward more plausible prediction directions. For instance, in personalized search scenarios, top-clicked items can be used as anchors to introduce collaborative filtering signals that reflect the aggregated behavior of many users, providing informative cues for capturing context-specific user intents.

Formally, let the current session of user  $u$  provide  $B$  anchor groups  $\mathcal{A}^u = \{A_1^u, \dots, A_B^u\}$  (e.g., top- $K$  purchased or clicked items related to the current user or query), where  $A_b^u = (v_{b,1}^{PA}, \dots, v_{b,m_b}^{PA})$  denotes the  $b$ -th group with  $m_b$  items. For the  $j$ -th anchor in group  $b$ , we compute the token embedding as:

$$\mathbf{z}_{b,j}^{PA} = \text{Proj}_{\text{shared}}(\phi_{\text{item}}(\mathbf{v}_{b,j}^{PA})) \in \mathbb{R}^d, \quad \mathbf{h}_{b,j}^{PA} = \mathbf{z}_{b,j}^{PA} + \mathbf{p}_j^{PA}. \quad (4)$$

Here,  $\mathbf{p}_j^{PA} \in \mathbb{R}^d$  represents the positional embedding for the  $j$ -th position within a group. To preserve group structure, we wrap each group with learnable boundary tokens  $\mathbf{e}_{\text{BOS}}, \mathbf{e}_{\text{EOS}} \in \mathbb{R}^d$  and get the final token sequence for each anchor group:

$$(\mathbf{e}_{\text{BOS}}, \mathbf{h}_{b,1}^{PA}, \dots, \mathbf{h}_{b,m_b}^{PA}, \mathbf{e}_{\text{EOS}}). \quad (5)$$

**Situational Descriptors (SD).** Situational context captures non-item information relevant to the ranking task. For the user  $u$  with its associated features, we compute the embedding:

$$\mathbf{z}^U = \text{Proj}_{\text{user}}(\phi_{\text{user}}(\mathbf{u})) \in \mathbb{R}^d, \quad \mathbf{h}^U = \mathbf{z}^U + \mathbf{p}_k^U, \quad (6)$$

where  $\mathbf{h}^U \in \mathbb{R}^d$  represents the projected user embedding and  $\mathbf{p}_k^U \in \mathbb{R}^d$  denotes the positional embedding for the user token at position  $k$  in the sequence.

Similarly, for the query  $q$  (which can be omitted in recommendation scenarios), we obtain:

$$\mathbf{z}^Q = \text{Proj}_{\text{query}}(\phi_{\text{query}}(\mathbf{q})) \in \mathbb{R}^d, \quad \mathbf{h}^Q = \mathbf{z}^Q + \mathbf{p}_k^Q, \quad (7)$$

where  $\mathbf{p}_k^Q \in \mathbb{R}^d$  is the corresponding positional embedding for the query token.

**Candidate Item Set (CIS, ranking mode only).** In the ranking stage, a key design choice in ranking is whether to process candidates in a *pointwise* or *setwise* manner.

- In the *pointwise* mode, each candidate item is concatenated with the user context and passed through the model independently. This design is computationally efficient and parallelizable across candidates, but it fails to directly compare between candidates, which is crucial for ranking.
- In the *full setwise* mode, the entire retrieved candidate set  $\mathcal{V}$  (often thousands of items) is encoded jointly. This allows the model to compare all candidates in a shared latent space, aligning better with the intrinsic nature of ranking. However, the computational cost is prohibitive, especially when combined with multi-step reasoning.

To strike a balance between efficiency and expressiveness, we adopt a *grouped setwise* strategy: the retrieved candidate set  $\mathcal{V}$  is randomly partitioned into smaller groups of size  $C$  (e.g., 12). Each group is processed independently in the ranking mode, enabling intra-group interaction among candidates while keeping the per-group reasoning cost tractable. Importantly, since grouping is randomized across training and inference, the model learns to be robust to arbitrary candidate organizations and achieves stable scoring performance. In this way, grouped setwise ranking thus strikes a practical

balance between pointwise ranking ( $C = 1$ ) and full-set ranking ( $C = |\mathcal{V}'|$ ), making it suitable for large-scale real-world personalized ranking systems deployment.

Formally, given a candidate set  $C^u = \{\mathbf{v}_1^{\text{CIS}}, \dots, \mathbf{v}_C^{\text{CIS}}\}$  containing  $C$  items to be ranked, we embed each candidate item as:

$$\mathbf{z}_i^{\text{CIS}} = \text{Proj}_{\text{cand}}(\phi_{\text{item}}(\mathbf{v}_i^{\text{CIS}})) \in \mathbb{R}^d. \quad (8)$$

To avoid spurious correlations between position and relevance labels, we deliberately exclude positional embeddings for candidate tokens:

$$\mathbf{h}_i^{\text{CIS}} = \mathbf{z}_i^{\text{CIS}}, \quad 1 \leq i \leq C. \quad (9)$$

This design choice ensures that the model evaluates each candidate purely based on its content rather than its position in the input sequence.

**Sequence Packing and Ordering.** Let  $\oplus$  denote the concatenation of token subsequences. The final input sequence to the backbone model is obtained by packing the Interaction Behaviors (IH), Preference Anchors (PA), Situational Descriptors (SD), and, in ranking mode only, the Candidate Item Set (CIS), following fixed ordering rules.

- *Retrieval Mode.* The final input token sequence is constructed as:

$$\mathcal{I}_{\text{retrieval}}^u = \underbrace{(\mathbf{h}_1^{\text{IH}}, \dots, \mathbf{h}_{n_u}^{\text{IH}})}_{\text{chronological IH}} \oplus \underbrace{\bigoplus_{b=1}^B (\mathbf{e}_{\text{BOS}}, \mathbf{h}_{b,1}^{\text{PA}}, \dots, \mathbf{h}_{b,m_b}^{\text{PA}}, \mathbf{e}_{\text{EOS}})}_{\text{PA groups ordered by business rule}} \oplus \underbrace{(\mathbf{h}^{\text{U}}, \mathbf{h}^{\text{Q}}, \dots)}_{\text{SD segment}}. \quad (10)$$

Here: (1) IH tokens are ordered by ascending interaction timestamp; (2) Each PA group is wrapped by BOS/EOS boundary tokens, with groups ordered according to predefined business rules; (3) SD tokens have no temporal ordering and are placed in a segment with distinct positional indices.

- *Ranking Mode.* We extend the retrieval-mode sequence by appending candidate item tokens:

$$\mathcal{I}_{\text{rank}}^u = \mathcal{I}_{\text{retrieval}}^u \oplus \underbrace{(\mathbf{h}_1^{\text{CIS}}, \dots, \mathbf{h}_C^{\text{CIS}})}_{\text{no positional embedding}}. \quad (11)$$

CIS tokens are appended without positional encodings to prevent the model from learning spurious correlations between sequence positions and relevance labels during training.

### 3.3. Backbone Architecture

Our backbone consumes the packed token sequence from Sec. 3.2 in a *unified* manner for both retrieval and ranking. Below, we first describe the sequence encoder and then the block-wise multi-step reasoning that exposes intermediate latent representations for multi-task progressive supervision learning discussed in Sec. 3.4.

#### 3.3.1. Transformer-Based Sequential Encoding

Let  $\mathcal{I} = [\mathbf{h}_1; \dots; \mathbf{h}_N]$  denote the final input tokens constructed from Sec. 3.2, where  $N$  is the length of the total input sequence. We adopt an  $L$ -layer bi-directional Transformer (Vaswani et al., 2017) with pre-normalization for the backbone architecture. Formally, let  $\mathbf{H}^l = [\mathbf{h}_1^l; \dots; \mathbf{h}_N^l] \in \mathbb{R}^{N \times d}$  denote the hidden states at later  $l$ , with initial input  $\mathbf{H}^{(0)} = \mathcal{I}$ , the  $l$ -th layer ( $1 \leq l \leq L$ ) is

$$\begin{aligned} \mathbf{H}_{\text{attn}}^l &= \mathbf{H}^{l-1} + \text{MHSA}\left(\text{LN}\left(\mathbf{H}^{l-1}\right)\right), \\ \mathbf{H}^l &= \mathbf{H}_{\text{attn}}^l + \text{FFN}\left(\text{LN}\left(\mathbf{H}_{\text{attn}}^l\right)\right), \end{aligned}$$

where  $\text{MHSA}(\cdot)$  is multi-head self-attention with *bi-directional* attention and  $\text{FFN}(\cdot)$  is a position-wise feed-forward network;  $\text{LN}(\cdot)$  denotes layer normalization (grouped LN can be used in practice). We prefer bi-directional attention because the personalized ranking task is non-autoregressive at both training and serving modes, and allowing tokens to condition on the full context yields superior performance in practice. The final encoder output  $\mathbf{H}^L = [\mathbf{h}_1^L; \dots; \mathbf{h}_N^L]$  serves as the foundation for subsequent reasoning.

### 3.3.2. Block-Wise Multi-Step Reasoning

Inspired by recent work on reasoning-enhanced recommendation (e.g., ReaRec (Tang et al., 2025)), we also employ a multi-step refinement process. Different from prior approaches that recycle a single hidden state across iterations, we design a *block-wise reasoning* mechanism, where a set of hidden states is iteratively transferred across steps. The motivation arises from the limitation of single-unit reasoning, which may overly compress signals due to the narrow capacity of its transmission medium, thereby losing important information. In contrast, block-wise reasoning offers adjustable reasoning bandwidth, providing greater flexibility and achieving a better balance between information compression and retention. Moreover, by grouping tokens into reasoning blocks, the model is encouraged to assign specialized roles to different tokens, leading to more structured and effective refinement of representations.

Formally, let  $M$  denotes the block size which is task-dependent, and let  $\mathbf{B}_k \in \mathbb{R}^{M \times d}$  denote the  $k$ -th ( $k \in \mathbb{Z}$ ) reasoning block. The initial block  $\mathbf{B}_0$  is constructed from the final encoder output  $\mathbf{H}^L$  as

$$\mathbf{B}_0 = \mathbf{H}^L[N - M + 1 : N] \in \mathbb{R}^{M \times d}.$$

For step  $k \geq 1$ , we first extract the block from the previous reasoning step output:

$$\mathbf{B}_k = \mathbf{H}_{k-1}^L[N + (k-2)M + 1 : N + (k-1)M] \in \mathbb{R}^{M \times d}.$$

To distinguish different reasoning steps, following the design of ReaRec Tang et al. (2025), we introduce Reasoning Position Embeddings (RPE).

Let  $\mathbf{E}_{\text{RPE}} \in \mathbb{R}^{K \times d}$  be a learnable embedding matrix where  $K$  is the maximum number of reasoning steps. We define:

$$\begin{aligned} \tilde{\mathbf{B}}_0 &= \mathbf{B}_0, \\ \tilde{\mathbf{B}}_k &= \mathbf{B}_k + \mathbf{1}_M \otimes \mathbf{E}_{\text{RPE}}[k, :] \quad \text{for } k \geq 1, \end{aligned}$$

where  $\mathbf{1}_M \in \mathbb{R}^M$  is a vector of ones and  $\otimes$  denotes the outer product. At each step  $k$ , we concatenate the base sequence  $\mathcal{I}$  with all previous enhanced blocks  $\tilde{\mathbf{B}}_{<k}$  and the current block  $\tilde{\mathbf{B}}_k$ , then pass them through the Transformer backbone with a block-wise causal mask:

$$[\mathcal{I}; \tilde{\mathbf{B}}_{<k}; \tilde{\mathbf{B}}_k] \xrightarrow{\mathcal{F}_\theta(\cdot; \mathcal{M}_k)} \mathbf{H}_k^L \in \mathbb{R}^{(N+kM) \times d}, \quad (12)$$

where  $\mathcal{F}_\theta$  is the Transformer update and  $\mathcal{M}_k$  is the mask mechanism that enforces information flow constraints. As shown in Figure 4(a), we adopt a causal block-wise mask  $\mathcal{M}_k$ , which enforces that (1)

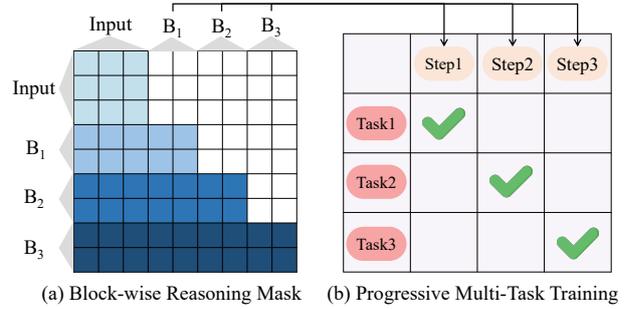


Figure 4 | Block-wise reasoning mask and progressive multi-task training. (a) Causal attention mask enables reasoning blocks to attend to input and previous blocks. (b) Progressive training assigns tasks of increasing complexity to successive reasoning steps to provide effective process supervision.

block tokens  $\tilde{\mathbf{B}}_k$  can attend to all base tokens  $\mathcal{I}$  and all historical blocks  $\tilde{\mathbf{B}}_{<k}$ ; (2) tokens of current reasoning block cannot attend to future reasoning block tokens. Iteratively applying this procedure yields progressively refined reasoning states  $\tilde{\mathbf{B}}_1, \tilde{\mathbf{B}}_2, \dots, \tilde{\mathbf{B}}_K$ .

The above formulation is generic, and we adopt task-specific strategies for block size  $M$  in industrial cascaded systems:

- *Retrieval Mode*. Here the block size  $M$  is set equal to the length of the situation descriptor (SD), which provides enough capacity for compact yet expressive refinement. Since retrieval serves as a coarse filtering stage, the objective is to strengthen and balance scenario-specific goals. Taking personalized search as an example, we designate the user ( $\mathbf{h}_U$ ) and query ( $\mathbf{h}_Q$ ) tokens as aggregation blocks:  $\mathbf{h}_U$  encodes personalization signals while  $\mathbf{h}_Q$  captures relevance, and iterative reasoning over these tokens allows the model to jointly reinforce both dimensions during matching.
- *Ranking Mode*. Here the objective is to distinguish subtle differences within a constructed candidate set, and the candidate item set (CIS) is fully visible to the model. We therefore set  $M = C$ , where each block corresponds to all candidate item tokens. The final block  $\tilde{\mathbf{B}}_K$  thus contains the refined representations used for ranking. To further prevent overfitting to candidate order and encourage robust set-wise reasoning, we apply randomized candidate grouping during training, which ensures the model learns reasoning abilities invariant to specific candidate arrangements within each group.

### 3.4. Progressive Multi-Task Training

Building upon the block-wise multi-step reasoning framework, we obtain a sequence of intermediate block representations  $\{\mathbf{B}_k\}_{k=1}^K$ , where  $\mathbf{B}_k \in \mathbb{R}^{M \times d}$  denotes the  $k$ -th reasoning block containing  $M$  token states of dimension  $d$ . To provide effective process supervisory signals for the multi-step reasoning trajectory, we introduce a progressive multi-task training paradigm that implements curriculum learning (Bengio et al., 2009) through gradually increasing task complexity.

We define  $K$  learning objectives  $\mathcal{T} = \{\tau_1, \tau_2, \dots, \tau_K\}$  arranged in a progressive curriculum from fundamental to advanced tasks. For example, in e-commerce domains, this follows the natural user engagement progression: exposure  $\rightarrow$  click  $\rightarrow$  purchase, where each successive task requires increasingly sophisticated understanding of user preferences and behavioral patterns. Each reasoning step  $k$  is assigned to optimize a single task  $\tau_k$ , creating a step-by-step learning trajectory where the model first masters basic recognition capabilities before advancing to complex preference modeling. This progressive assignment provides rich process supervision at each reasoning stage, enabling more effective learning of intermediate representations and ensuring that each reasoning block develops specialized capabilities while building upon the foundation established by previous steps.

#### 3.4.1. Retrieval Mode

In the retrieval stage, we extract user representations from the sequence of reasoning blocks and optimize them using a combination of calibrated probability estimation and bidirectional contrastive learning objectives.

Specifically, given the sequence of reasoning blocks  $\{\mathbf{B}_k\}_{k=1}^K$  where  $\mathbf{B}_k \in \mathbb{R}^{M \times d}$ , we extract a step-specific user representation from each block via layer normalization followed by mean pooling:

$$\mathbf{r}_k = \text{Mean}(\text{LN}(\mathbf{B}_k)) \in \mathbb{R}^d, \quad k \in \{1, 2, \dots, K\}.$$

For each training instance, we construct a candidate pool  $\Omega$  containing items with diverse behavioral labels across the progressive task sequence. Let  $\mathbf{z}_v \in \mathbb{R}^d$  denote the item embedding for candidate

$v \in \Omega$ . For the task  $\tau_k$  assigned to step  $k$  and candidate  $v \in \Omega_{\tau_k}$ , let  $y_v^k \in \{0, 1\}$  denote the corresponding behavioral label. We partition the candidate pool into positive and negative sets:

$$\Omega_{\tau_k}^+ = \{v \in \Omega_{\tau_k} : y_v^k = 1\}, \quad \Omega_{\tau_k}^- = \{v \in \Omega_{\tau_k} : y_v^k = 0\}.$$

We employ two complementary learning objectives that address different granularities of representation optimization:

(i) **Binary Cross-Entropy Loss (BCE)** operates at the point-wise level, providing calibrated probability estimates for individual user-item pairs:

$$\mathcal{L}_k^{\text{BCE}} = \sum_{v \in \Omega_{\tau_k}^+} -\log \sigma(\langle \mathbf{r}_k, \mathbf{z}_v \rangle) + \sum_{v \in \Omega_{\tau_k}^-} -\log (1 - \sigma(\langle \mathbf{r}_k, \mathbf{z}_v \rangle)), \quad (13)$$

where  $\sigma(\cdot)$  is the sigmoid function.

(ii) **Bidirectional Contrastive Learning (BCL)** operates at the batch level, enabling global contrastive reasoning across the in-batch samples. Drawing inspiration from CLIP (Radford et al., 2021), we propose BCL optimization objective that consists of two symmetric components:

**User-to-Item (U2I) Contrastive Learning** enables each user representation  $\mathbf{r}_k$  to holistically distinguish positive items from negative candidates, which is formulated as:

$$\mathcal{L}_k^{\text{U2I}} = \sum_{v \in \Omega_{\tau_k}^+} -\log \frac{\exp(\langle \mathbf{r}_k, \mathbf{z}_v \rangle / \eta)}{\sum_{v^+ \in \Omega_{\tau_k}^+} \exp(\langle \mathbf{r}_k, \mathbf{z}_{v^+} \rangle / \eta) + \sum_{v^- \in \Omega_{\tau_k}^-} \exp(\langle \mathbf{r}_k, \mathbf{z}_{v^-} \rangle / \eta)}, \quad (14)$$

where  $\eta > 0$  is the temperature parameter.

**Item-to-User (I2U) Contrastive Learning** enables each positive item to identify its corresponding user from all in-batch users. Let  $\mathcal{R}_k = \{\mathbf{r}_k^{(i)}\}_{i=1}^B$  denote the set of step- $k$  user representations within the current training batch, where  $B$  is the batch size. Then, the I2U objective is:

$$\mathcal{L}_k^{\text{I2U}} = \sum_{v \in \Omega_{\tau_k}^+} -\log \frac{\exp(\langle \mathbf{r}_k, \mathbf{z}_v \rangle / \eta)}{\sum_{\mathbf{r}' \in \mathcal{R}_k} \exp(\langle \mathbf{r}', \mathbf{z}_v \rangle / \eta)}. \quad (15)$$

The complete BCL objective for step  $k$  combines both symmetric components:

$$\mathcal{L}_k^{\text{BCL}} = \mathcal{L}_k^{\text{U2I}} + \mathcal{L}_k^{\text{I2U}}. \quad (16)$$

And the overall retrieval loss aggregates the objectives across all reasoning steps, where each step  $k$  optimizes for its assigned task  $\tau_k$ :

$$\mathcal{L}^{\text{retrieval}} = \sum_{k=1}^K (\mathcal{L}_k^{\text{BCE}} + \mathcal{L}_k^{\text{BCL}}). \quad (17)$$

This progressive training strategy provides comprehensive process supervision, where each reasoning step receives dedicated supervisory signals tailored to its assigned task complexity, achieving more effective learning of multi-step reasoning capabilities.

### 3.4.2. Ranking Mode

In the ranking stage, the block size equals the candidate group size ( $M = C$ ). Each reasoning block  $\mathbf{B}_k \in \mathbb{R}^{C \times d}$  contains  $C$  hidden states  $\{\mathbf{h}_{i,k}\}_{i=1}^C$ , where  $\mathbf{h}_{i,k} \in \mathbb{R}^d$  represents the hidden state for the  $i$ -th

candidate at reasoning step  $k$ . For the task  $\tau_k$  assigned to reasoning step  $k$ , we compute candidate-wise logits through a task-specific scoring network:

$$s_{i,k} = \text{MLP}_{\tau_k}(\mathbf{h}_{i,k}) \in \mathbb{R}, \quad i \in \{1, \dots, C\},$$

where  $\text{MLP}_{\tau_k}$  denotes the Multi-Layer Perceptron (MLP) for task  $\tau_k$ .

We employ two complementary learning objectives:

(i) **Binary Cross-Entropy Loss (BCE)** provides point-wise probability calibration for individual candidates:

$$\mathcal{L}_k^{\text{BCE}} = \sum_{i=1}^C [-y_i^{\tau_k} \log \sigma(s_{i,k}) - (1 - y_i^{\tau_k}) \log(1 - \sigma(s_{i,k}))], \quad (18)$$

where  $y_i^{\tau_k} \in \{0, 1\}$  denotes the behavioral label of candidate  $i$  for task  $\tau_k$ .

(ii) **Set Contrastive Learning (SCL)** operates at the set-wise level, enabling each positive candidate to distinguish itself from negative candidates within the set:

$$\mathcal{L}_k^{\text{SCL}} = \sum_{i: y_i^{\tau_k}=1} -\log \frac{\exp(s_{i,k}/\eta)}{\sum_{j=1}^C \exp(s_{j,k}/\eta)}, \quad (19)$$

where the summation is over all positive candidates, and each positive candidate competes against all candidates in the group for ranking position.

Then, the overall ranking loss combines both objectives across all reasoning steps:

$$\mathcal{L}^{\text{ranking}} = \sum_{k=1}^K (\mathcal{L}_k^{\text{BCE}} + \mathcal{L}_k^{\text{SCL}}). \quad (20)$$

This curriculum-based multi-task training provides dense supervision signals throughout the reasoning process, enabling the model to learn more effective step-by-step reasoning patterns while ensuring each reasoning block develops appropriate capabilities for its designated task complexity level.

### 3.5. Time Complexity Analysis

In this section, we analyze the complexity of our block-wise reasoning framework under retrieval and ranking modes.

For the backbone encoder (without reasoning), the time complexity of each Transformer layer is

$$O(N^2d + Nd^2),$$

where  $N$  is the base sequence length,  $d$  is the hidden dimension, and  $L$  is the number of layers. Thus the total cost is  $O(L(N^2d + Nd^2))$ .

For the reasoning phase, we employ the *KV Caching* technique to reuse historical key-value pairs, so that each new reasoning step only incurs attention between the  $M$  new block tokens and the cached tokens. Specifically, at reasoning step  $k$ , we need to: (i) Compute Q, K, V for  $M$  new block tokens:  $O(Md^2)$ ; (ii) Calculate attention between  $M$  new tokens and all  $N + (k-1)M$  cached tokens:  $O(M(N + kM)d)$ ; (iii) Apply output projection:  $O(Md^2)$ . Therefore, the complexity per layer at step  $k$  is:

$$O(M(N + kM)d + Md^2).$$

Aggregating over  $K$  reasoning steps and  $L$  layers, the total additional reasoning cost is

$$O(LKM(Nd + MKd + d^2)).$$

## 4. Offline Experiments

To validate the effectiveness of OnePiece, we first conduct extensive offline experiments based on 30-day log from Shopee, an e-commercial platform serving billion-scale users across Southeast Asia and Latin America. The full dataset statistics are summarized in Table 1, and the detailed construction process of training samples is described in Appendix A. Furthermore, we perform comprehensive ablation studies and other detailed analyses to reveal the advantages of OnePiece.

Table 1 | Dataset Statistics of Shopee E-commerce Platform from June 11 to July 10, 2025. The dataset includes multi-behavior user interaction. Abbreviations: M = Million ( $10^6$ ), B = Billion ( $10^9$ ).

#User	#Item	#Query	#Impression	#Click	#Add-to-Cart	#Order
10M	93M	12M	0.24B	60M	12M	6M

### 4.1. Experimental Settings

#### 4.1.1. Evaluation Protocol

We adopt a streaming-style training and evaluation protocol to closely mimic real-world deployment. For training, we employ an incremental scheme where the model is trained day by day, with samples within each day randomly shuffled to ensure robustness. For offline evaluation, we follow a strict chronological order: the checkpoint trained on day- $t$  is evaluated on the unseen data from day- $(t + 1)$ . Unless otherwise specified, all experiments are conducted on the same 30 consecutive days of data to ensure comparability across methods.

#### 4.1.2. Evaluation Metrics

For the retrieval stage, we are primarily concerned with how many clicked items are successfully recalled. Thus, we report Recall@100 and Recall@500, denoted as R@100 and R@500. For the ranking stage, we evaluate the model under three types of user feedback: click (C-), add-to-cart (A-), and order (O-). For each type, we report both AUC and GAUC, denoted as C-AUC/C-GAUC, A-AUC/A-GAUC, and O-AUC/O-GAUC, respectively.

#### 4.1.3. Baseline Methods

We compare OnePiece against several representative baselines:

**DLRM** (Production baseline in Shopee): Shopee’s DLRM baseline is a well-optimized hybrid model including multiple SOTA components. For *retrieval mode*, it adopts a two-tower design inspired by DSSM (Huang et al., 2013), where the query and user history are encoded separately. To enhance relevance modeling, it employs a DIN-like (Zhou et al., 2018) structure and zero-attention (Ai et al., 2019). A lightweight text CNN is used to extract keyword features, while DCNv2 (Wang et al., 2021) captures high-order cross features. Sequential features are aggregated by mean pooling, concatenated with other features, and fused via an MLP. For *ranking mode*, the DLRM model is different because we can allow incorporating the candidate item into a single tower with user features. The backbone is ResFlow (Fu et al., 2024), combined with DIN-like target attention and cross-attention across sequential behaviors. DCNv2 is again used for higher-order interactions, followed by MLP fusion. A SENet (Hu et al., 2017) module further supports adaptive feature selection for different tasks.

**HSTU** (Zhai et al., 2024): HSTU is a representative generative recommendation framework proposed by Meta. For a fair comparison, we align its parameter size with OnePiece and adopt the

Table 2 | Performance comparison of different models on both retrieval and ranking tasks with 30 days of training data. For retrieval, we report Recall@100 (R@100) and Recall@500 (R@500) evaluated on clicked samples. For ranking, we report AUC and GAUC under three feedback types: click (C-), add-to-cart (A-), and order (O-). The best results are highlighted in **bold**.

Model	Retrieval Mode		Ranking Mode					
	R@100	R@500	C-AUC	C-GAUC	A-AUC	A-GAUC	O-AUC	O-GAUC
DLRM	0.458	0.679	0.856	0.851	0.893	0.843	0.931	0.854
HSTU	0.443	0.658	0.833	0.829	0.878	0.827	0.913	0.839
HSTU+PA	0.472	0.680	0.855	0.852	0.901	0.848	0.926	0.849
ReaRec	0.452	0.674	0.843	0.838	0.882	0.834	0.919	0.843
ReaRec+PA	0.485	0.701	0.862	0.863	0.908	0.851	0.927	0.851
<b>OnePiece</b>	<b>0.517</b>	<b>0.731</b>	<b>0.911</b>	<b>0.909</b>	<b>0.952</b>	<b>0.897</b>	<b>0.963</b>	<b>0.886</b>

same side-information fusion strategy. Since the original HSTU only considers Interaction History (IH) and Situational Descriptors (SD), we additionally evaluate a variant HSTU+PA where Preference Anchors (PA) are introduced into its input sequence, consistent with OnePiece’s context engineering.

**ReaRec** (Tang et al., 2025): ReaRec is a representative reasoning-enhanced recommendation model that formulates user representation modeling as multi-step reasoning over item sequences. The vanilla ReaRec architecture only supports retrieval tasks with the user iteration history sequence as input. To ensure experimental consistency, we adapt its backbone and feature inputs to match OnePiece. For ranking, we adapt ReaRec by introducing the candidate item into the input sequence and applying a target-aware attention mask (following the design in HSTU), where sequence tokens can attend to the candidate but candidates remain mutually invisible. Similar to HSTU, we also evaluate a ReaRec+PA variant that augments IH and SD with Preference Anchors, providing a stronger context prior for reasoning.

## 4.2. Overall Performance

Table 2 summarizes the performance of different models on both retrieval and ranking tasks. Several observations can be drawn as follows:

First, **DLRM remains a very strong baseline**. Compared with original HSTU and ReaRec, DLRM achieves better performance on most evaluated metrics because it fully exploits rich feature interactions, target attention, and multiple sequential features. In contrast, the original HSTU and ReaRec underperform DLRM as they only rely on interaction behavior sequences as input.

Second, **the preference anchor (PA) mechanism brings consistent benefits independent of the backbone**. Both HSTU+PA and ReaRec+PA outperform their vanilla counterparts, confirming that enriching user history with auxiliary preference signals provides complementary information. Between the two, ReaRec demonstrates higher robustness due to its Transformer backbone with bidirectional attention and reasoning capability, which enables it to better leverage the anchor information than HSTU.

Finally, **OnePiece achieves the best overall results**. Compared with the strongest ReaRec+PA baseline, OnePiece further improves Recall@100 from 0.485 to 0.517 and C-AUC from 0.862 to 0.911. These consistent gains can be attributed to its novel block-wise latent reasoning and the progressive multi-task training strategy, which enable finer-grained preference refinement through multi-step reasoning. Together, these advances validate the design of OnePiece as a more powerful

Table 3 | Ablation studies on context engineering of OnePiece. PA( $L$ ) denotes appending a preference-anchor sequence with maximum length  $L$ . SD denotes situational descriptors.

Version	Model	Retrieval		Ranking					
		R@100	R@500	C-AUC	C-GAUC	A-AUC	A-GAUC	O-AUC	O-GAUC
V1	IH(ID)	0.407	0.646	0.802	0.802	0.860	0.819	0.908	0.835
V2	IH(ID+Side Info)	0.428	0.657	0.846	0.844	0.871	0.839	0.918	0.845
V3	V2+PA(10)	0.459	0.677	0.879	0.876	0.923	0.863	0.940	0.861
V4	V2+PA(20)	0.467	0.686	0.885	0.886	0.929	0.869	0.946	0.866
V5	V2+PA(30)	0.475	0.689	0.892	0.890	0.936	0.874	0.949	0.871
V6	V2+PA(60)	0.491	0.707	0.901	0.900	0.945	0.886	0.956	0.880
V7	V2+PA(90)	0.504	0.719	0.908	0.905	0.951	0.896	0.962	0.885
V8	V7+SD	<b>0.517</b>	<b>0.731</b>	<b>0.911</b>	<b>0.909</b>	<b>0.952</b>	<b>0.897</b>	<b>0.963</b>	<b>0.886</b>

reasoning-enhanced recommendation framework capable of unifying retrieval and ranking.

### 4.3. Ablation Study

In this section, we provide detailed ablation studies to evaluate the contribution of each design in OnePiece through a step-by-step analysis.

#### 4.3.1. Context Engineering Ablation

To validate the effectiveness of our context engineering design, we conduct a step-by-step ablation study, progressively adding Interaction History (IH), Preference Anchors (PA), and Situational Descriptors (SD) into the model input. The experiment results for both retrieval and ranking are summarized in Table 3.

We begin with a minimal baseline (V1), where only user interaction sequences composed of raw item IDs are utilized. This configuration, implemented with a two-layer Transformer and bidirectional attention, yields the lowest performance across both retrieval and ranking tasks. Furthermore, introducing side information for each item (V2) leads to a clear improvement, highlighting the importance of more item features beyond raw IDs.

Building on this, we incorporate Preference Anchors (PA) by appending users’ query-associated item sequences comprising top- $k$  purchases, top- $k$  clicks, and top- $k$  impressions. From V3 to V7, we gradually increase the sequence length from 10 to 90, progressively extending the combined behavioral sequence. The results demonstrate a clear **scaling effect of PA**: extending the auxiliary item sequence enriches query-specific context, enabling the model to capture more fine-grained user intent and leading to steadily higher performance. In particular, PA provides query-dependent signals that are absent in plain IH, allowing the model to differentiate user preferences under different queries rather than relying solely on global interaction history or generic query descriptors.

Finally, incorporating Situational Descriptors (V8) yields the best overall results, especially on retrieval. Compared with V7, Recall@100 improves from 0.504 to 0.517, as the additional user- and query-side descriptors provide stronger contextual grounding that helps retrieve a broader set of relevant items. By contrast, the ranking gains are marginal, since IH already supplies rich personalization signals and PA captures detailed query-specific preferences. Given that ranking primarily focuses on fine-grained comparisons among candidate items where query relevance is high, SD serves as a weaker anchor in this stage.

Table 4 | Effect of different training strategies on retrieval performance. Here, multi-task refers to jointly optimizing impression and click prediction losses.

Version	Training Strategy	R@100	R@500
V1	Causal Mask	0.464	0.671
V2	Bi-Directional	0.470	0.676
V3	V2 + 1-Step Reasoning, Click Task on Last Step	0.490	0.708
V4	V2 + 1-Step Reasoning, Multi-Task on Last Step	0.495	0.714
V5	V2 + 2-Step Reasoning, Multi-Task on Last Step	0.510	0.726
V6	V2 + 2-Step Reasoning, Progressive Multi-Task	<b>0.517</b>	<b>0.731</b>

Table 5 | Effect of different training strategies on ranking performance. Here, multi-task refers to jointly optimizing impression, click, and order prediction losses.

Version	Training Strategy	C-AUC	C-GAUC	A-AUC	A-GAUC	O-AUC	O-GAUC
V1	Causal Mask	0.839	0.836	0.876	0.830	0.911	0.838
V2	Bi-Directional, CIS Inter-Invisible	0.860	0.859	0.903	0.848	0.920	0.847
V3	Bi-Directional, CIS Inter-Visible	0.881	0.879	0.918	0.857	0.937	0.854
V4	V3 + 1-Step Reasoning, Multi-Task on Last Step	0.890	0.889	0.931	0.871	0.946	0.867
V5	V3 + 2-Step Reasoning, Multi-Task on Last Step	0.893	0.894	0.936	0.876	0.948	0.869
V6	V3 + 3-Step Reasoning, Multi-Task on Last Step	0.906	0.902	0.946	0.889	0.957	0.881
V7	V3 + 3-Step Reasoning, Progressive Multi-Task	<b>0.911</b>	<b>0.909</b>	<b>0.952</b>	<b>0.897</b>	<b>0.963</b>	<b>0.886</b>

Overall, these ablation studies demonstrate that **structured context engineering is highly effective**: IH captures long-term personalization, PA provides scalable and query-specific anchors, and SD contributes stable contextual grounding, particularly benefiting retrieval. Together, they enrich user–query representations in complementary ways and enable OnePiece to achieve consistent improvements across retrieval and ranking.

#### 4.3.2. Training Strategy Ablation

We systematically examine the impact of different training strategies on both retrieval and ranking tasks to validate our progressive multi-task learning framework. The results are demonstrated in Table 4 and Table 5, respectively.

Starting from a causal attention baseline (V1), bidirectional attention (V2) effectively unlocks substantial gains across both tasks, improving retrieval R@100 from 0.464 to 0.470 and ranking C-AUC from 0.839 to 0.860. This validates our architecture design leveraging the bidirectional attention nature of recommendation tasks, where full context attention provides more comprehensive representation information. For ranking tasks, enabling candidate inter-visibility (V3) delivers another major boost, with C-AUC jumping from 0.860 to 0.881. This substantial improvement confirms that our candidate-aware framework enables the rich comparative reasoning essential for accurate ranking, allowing candidates to attend to each other rather than being processed in isolation.

The introduction of our block-wise reasoning mechanism demonstrates consistent and cumulative performance gains across both tasks. In retrieval, extending from direction decision-making to single-step reasoning with click prediction (V3) achieves notable improvements (R@100: 0.470 to 0.490), while multi-task learning on the final step (V4) provides additional gains (R@100: 0.495). Increasing reasoning depth yields expected progressive benefits: two-step reasoning (V5) improves retrieval R@100 to 0.510, while three-step reasoning in ranking (V6) achieves C-AUC of 0.906. These results establish that our reasoning framework enables increasingly sophisticated preference modeling, where each additional reasoning step meaningfully contributes to performance by capturing more nuanced

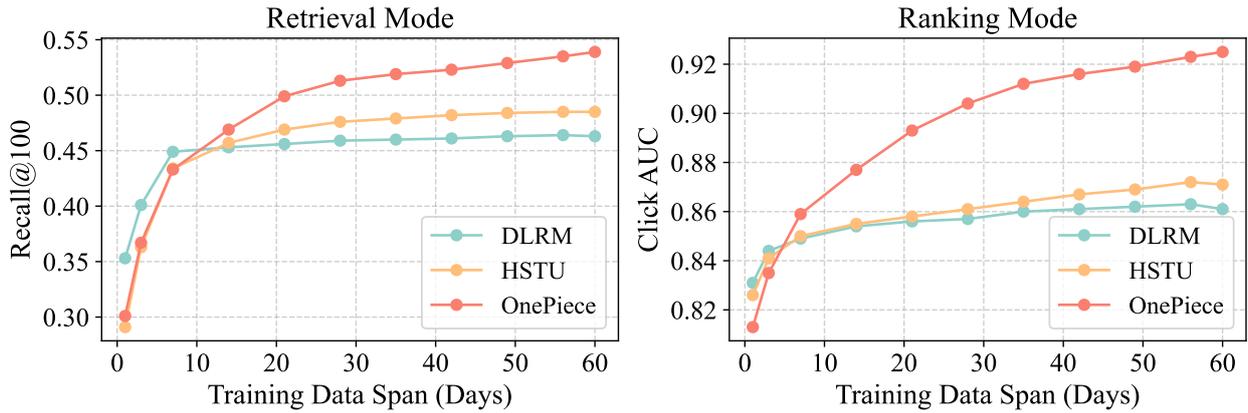


Figure 5 | Training convergence curves of different models on retrieval and ranking tasks.

user behavioral patterns.

Our progressive multi-task training consistently outperforms single-embedding multi-task learning at final step across both tasks. The key advantage lies in distributing different tasks across multiple reasoning steps rather than concentrating all supervision on a single final embedding. Progressive training surpasses the single-embedding approach with notable margins, improving R@100 from 0.510 (V5) to 0.517 (V6) in retrieval and C-AUC from 0.906 (V6) to 0.911 (V7) in ranking. This design prevents gradient conflicts that arise when multiple tasks compete for optimization at the one representation. Instead, each reasoning step serves as a specialized read-out token that helps the model extract task-specific information, effectively decoupling gradient flows across different objectives. The results also reveal task-specific optimization patterns: retrieval benefits from focused two-step reasoning aligned with the impression-click hierarchy, while ranking achieves optimal performance with three-step reasoning capturing the full conversion funnel, demonstrating how our progressive framework naturally adapts to different task complexities.

#### 4.4. Scaling Analysis

In this section, we provide more in-depth analyses to better understand the effectiveness of OnePiece.

##### 4.4.1. Training Data Scaling

To analyze how performance scales with increasing training data, we report Recall@100 for retrieval and AUC for ranking over different training spans (*i.e.*, the number of consecutive log days used for training), covering up to 60 days of data. To ensure a fair comparison, we carefully match the model size (*i.e.*, the total number of trainable parameters) across DLRM, HSTU, and our proposed OnePiece. As shown in Figure 5, OnePiece already surpasses both baselines after only 7–10 days of training, demonstrating superior data efficiency attributable to its context-aware design and multi-step reasoning architecture. With longer training spans, DLRM and HSTU quickly converge to a plateau, whereas OnePiece continues to improve with a widening margin. By day 60, the performance gap between OnePiece and baseline methods (DLRM and HSTU) becomes pronounced, with the baseline models exhibiting convergence behavior while OnePiece continues to demonstrate improvement potential. This suggests that OnePiece has not yet reached full convergence and can achieve further performance gains with additional training data, indicating superior scaling capabilities. These results highlight OnePiece’s stronger modeling capacity and its ability to effectively exploit richer behavioral supervision from extended time horizons. The trend is consistent across both retrieval

Table 6 | Effect of block size on ranking performance of OnePiece. The training data span is 60 days for this experiment.

Block Size	C-AUC	C-GAUC	A-AUC	A-GAUC	O-AUC	O-GAUC
$M = C = 1$	0.885	0.881	0.923	0.861	0.947	0.871
$M = C = 4$	0.913	0.911	0.951	0.896	0.961	0.885
$M = C = 8$	0.920	0.918	0.956	0.899	0.964	0.887
$M = C = 12$	0.927	0.923	0.958	0.903	0.969	0.893

and ranking tasks, and the training curves of OnePiece exhibit smooth and stable growth without noticeable fluctuations, indicating robust optimization under progressive multi-task supervision. Overall, OnePiece not only achieves higher sample efficiency but also scales more effectively as additional training data becomes available.

#### 4.4.2. Reasoning Scaling

We investigate the impact of reasoning block size  $M$  (i.e., the number of candidate items considered in ranking mode<sup>1</sup>) on the ranking performance of OnePiece. As shown in Table 6, increasing  $M$  from 1 to 12 yields consistent improvements across all evaluated metrics. The largest performance gain appears when scaling from  $M = 1$  to  $M = 4$ , since pointwise modeling ( $M = 1$ ) lacks cross-sample comparisons, while grouping candidates into blocks enables the reasoning mechanism to contrast preferences more effectively. As block size continues to increase, the improvements become smaller yet remain positive, indicating diminishing returns. This effect likely arises because overly large blocks overload the reasoning medium with redundant information, saturating its representational capacity. Overall, these findings reveal a trade-off between expanding reasoning bandwidth and avoiding information redundancy, underscoring the importance of selecting an appropriate block size to maximize the effectiveness of block-wise reasoning.

## 5. Online A/B Testing

In this section, we conduct large-scale online A/B testing on Shopee’s production search system to further evaluate the real-world effectiveness of OnePiece.

### 5.1. Experimental Settings

#### 5.1.1. Online Inference Details

**Retrieval Stage.** For retrieval, we conduct offline training to generate vector representations for the entire item pool, constructing an Approximate Nearest Neighbor (ANN) (Indyk and Motwani, 1998) index to support efficient online retrieval. Specifically, we employ the Hierarchical Navigable Small World (HNSW) (Malkov and Yashunin, 2018) algorithm.

**Ranking Stage.**<sup>2</sup> For ranking, we adopt the following score fusion strategy to integrate outputs from

<sup>1</sup>Note that we only conduct block size scaling experiments in the ranking stage. In retrieval, the candidate items are not fed into the recommendation model, and each block contains only two global tokens (a user token capturing personalization and a query token capturing relevance). This design leaves little room for block-size ablation in retrieval mode.

<sup>2</sup>Due to resource constraints in online deployment, the production version of OnePiece ranking model is a downgraded variant where the block size is reduced to  $M=1$  and the number of item side-information features is significantly reduced.

different tasks:

$$p_{\text{final}} = \alpha \cdot p_{\text{ctr}}^a \cdot p_{\text{ctcvr}}^b + \beta \cdot p_{\text{ctr}}^a \cdot p_{\text{ctcvr}}^b \cdot \text{price} + \gamma \cdot p_{\text{ctr}} \cdot \text{ecpm}, \quad (21)$$

where  $\alpha$ ,  $\beta$ , and  $\gamma$  are hyperparameters controlling the importance weights of respective components, enabling the balance between user experience and business revenue.  $p_{\text{ctr}}$  and  $p_{\text{ctcvr}}$  represent the click-through rate and click-to-conversion rate predicted by OnePiece’s final reasoning step, corresponding to the logits of the click and order tasks respectively. Parameters  $a$  and  $b$  modulate the influence of each task in the final ranking, while  $\text{price}$  denotes the item’s price information and  $\text{ecpm}$  represents the item’s advertising value component.

### 5.1.2. Evaluation Protocol

In industrial personalized ranking systems, candidate generation typically follows a cascade paradigm consisting of multiple stages (e.g., retrieval, pre-ranking, ranking). The retrieval stage usually generates a large set of candidate items through multiple parallel recall strategies, each capturing different facets of user preference to ensure diversified coverage, while the subsequent pre-ranking and ranking stages aim to produce the final ranked list with fine-grained preference estimation. As an initial online exploration, we allocate 10% of traffic for all A/B testing experiments under the following settings:

- For retrieval evaluation, we replace one of the parallel recall strategies with our proposed OnePiece retrieval model, specifically replacing the User-to-Item (U2I) recall route originally implemented in DLRM.
- For ranking evaluation, we substitute the DLRM model in the pre-ranking stage with our proposed OnePiece ranking model.

These designs allow us to isolate the contribution of OnePiece at different modes (i.e., retrieval and ranking) within the industrial cascade ranking systems.

### 5.1.3. Evaluation Metrics

We report the following business and user engagement indicators tracked in our online experiment:

- **GMV/UU** (Gross Merchandise Volume per Unique User): Average transaction value contributed per user.
- **GMV(99.5%)/UU**: GMV per user excluding the top 0.5% high-value orders, reflecting stable contributions from regular transactions.
- **AR/UU** (Advertising Revenue per Unique User): Average advertising revenue generated per unique user, reflecting conversion efficiency from ads exposure.
- **Order/UU**: Average number of orders per user, capturing transaction frequency.
- **Paid Order/UU**: Average number of successfully paid orders per user, counting only completed purchases without refund.
- **CTR** (Click-Through-Rate): The ratio of clicked impressions to total impressions, reflecting the effectiveness of ranked results in attracting user engagement.
- **CTCVR** (Click-to-Conversion Rate): The ratio of successful conversions to total clicks, measuring the effectiveness of transforming user engagement into completed transactions.
- **Buyer**: Proportion of unique users who placed at least one order.
- **Bad Query Rate**: The percentage of queries for which human evaluators judge the recommended content as irrelevant to the query, serving as an inverse measure of recommendation accuracy and user satisfaction.

Table 7 | Results of online A/B testing for the retrieval mode of OnePiece. Improvements are reported in relative percentage changes over the DLRM baseline.

GMV/UU	GMV(99.5%)/UU	Order/UU	Paid Order/UU	CTCVR	Buyer	Bad Query Rate
+1.08%	+0.91%	+0.71%	+0.98%	+0.66%	+0.41%	-0.17%

Table 8 | Results of online A/B testing for the ranking mode of OnePiece. Improvements are reported in relative percentage changes over the DLRM baseline.

GMV/UU	GMV(99.5%)/UU	AR/UU	Order/UU	Buyer	CTR	Bad Query Rate
+1.12%	+0.65%	+2.90%	+0.08%	+0.08%	+0.29%	+0.21%

## 5.2. Overall Performance

**Retrieval Mode.** Table 7 reports the detailed results of our online A/B testing for the retrieval mode of OnePiece. We observe consistent improvements across all key business metrics. GMV/UU increases by +1.08%, and even after removing the top 0.5% of high-value orders, GMV(99.5%)/UU still improves by +0.91%, suggesting that the gain is not driven solely by occasional big-ticket purchases but reflects stable contributions from regular transactions. Order/UU rises by +0.71%, while Paid Order/UU grows even faster (+0.98%) than Order/UU, indicating both higher conversion and a reduced refund rate. Buyer expands by +0.41%, meaning that more unique users were successfully converted. At the same time, CTCVR improves by +0.66%, reflecting a higher end-to-end conversion rate from exposure to conversion. Importantly, the Bad Query Rate decreases by 0.17%, which indicates better query relevance and improved user experience. Unlike prior personalized recall strategies that often boosted GMV at the expense of relevance, OnePiece achieves balanced improvements, simultaneously enhancing personalization, relevance, and overall transaction stability.

**Ranking Mode.** Table 8 summarizes the online A/B testing results of deploying OnePiece in the ranking stage. We observe consistent lifts across major business metrics: GMV/UU improves by +1.12%, and advertising revenue (AR/UU) shows a substantial boost of +2.9%. In contrast, Order/UU increases only marginally (+0.08%), which aligns with our design that translates order-related utility into GMV and advertising gains via the score fusion function (Eq. 21). Buyer also rises slightly by +0.08%, indicating broader coverage of converted users. On engagement metrics, CTR improves by +0.29%, suggesting stronger attractiveness of ranked results. Meanwhile, the Bad Query Rate increases by +0.21%, likely due to more advertising slots introducing items less relevant to user interests. However, this minor increase is outweighed by the substantial revenue gains. Overall, these results demonstrate that OnePiece ranking not only strengthens core business metrics but also achieves a practical trade-off between user experience and business objectives in large-scale industrial ranking systems.

## 5.3. Recall Coverage and Exclusive Contribution

To further evaluate the effectiveness of OnePiece in the retrieval stage, we analyze the overlap between the exposure items contributed by OnePiece and those from other recall routes (*i.e.*, recall coverage), in comparison with the traditional DLRM baseline. As shown in Table 9, OnePiece consistently achieves higher recall coverage across all routes. In sparse text recall, it improves STR1 coverage from 37.3% to 66.2% (+77.6%) and STR2 coverage from 31.3% to 64.4% (+105.8%). Similar gains are observed in graph-based Swing I2I (+32.6%) and keyword-popularity recall KPop (+23.5%), while semantic recall S2I also increases significantly from 47.6% to 67.8% (+42.4%). These consistent improvements

Table 9 | Comparison of overlap coverage between DLRM and OnePiece with respect to other recall strategies. For each recall route  $R$ , the recall coverage is computed as  $\text{Coverage}(R) = \frac{|\text{Exposure}_{\text{U2I}} \cap \text{Exposure}_R|}{|\text{Exposure}_R|}$ . Here,  $\text{Exposure}_{\text{U2I}}$  denotes the exposure set attributed to the U2I model (DLRM or OnePiece), and  $\text{Exposure}_R$  denotes the exposure set of another recall route. **STR1**: sparse text recall with user-input keywords; **STR2**: sparse text recall with rewritten keywords; **Swing I2I**: graph-based item-to-item personalized recall; **KPop**: popularity-based recall under keywords; **S2I**: semantic vector-to-item recall. Relative improvements of OnePiece over DLRM are shown in red.

Recall Route	STR1	STR2	Swing I2I	KPop	S2I
DLRM	37.3%	31.3%	57.9%	62.5%	47.6%
OnePiece	66.2% (+77.6%)	64.4% (+105.8%)	76.8% (+32.6%)	77.2% (+23.5%)	67.8% (+42.4%)

demonstrate that, compared to DLRM, OnePiece achieves substantial coverage gains over all other recall routes, revealing its strong potential to replace multiple specialized recall strategies with a single unified model. We believe that, with carefully designed context engineering, OnePiece can effectively balance personalization (Swing I2I), popularity (KPop), and relevance (STR1 and STR2), thereby moving closer to a “one model serves all purposes” paradigm for industrial-scale recall.

To further examine the independent value of OnePiece, we compare the exclusive contribution of OnePiece and DLRM in terms of impressions and clicks, as shown in Figure 6. OnePiece demonstrates a substantial increase in unique contribution: exclusive impression share rises from 3.6% to 9.9% (2.8x), while exclusive click share grows from 2.4% to 5.7% (2.4x). These results indicate that OnePiece is not only capable of covering the exposure of other recall routes but also contributes significantly more novel impressions and clicks that are not captured by traditional DLRM-based recall. In other words, the new recall route of OnePiece nearly doubles the independent value over traditional DLRM, showcasing its effectiveness in enhancing overall recall performance.

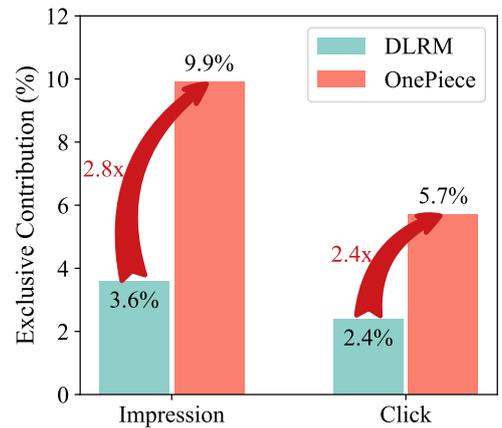


Figure 6 | Exclusive contribution of OnePiece in the retrieval stage.

#### 5.4. Efficiency Analysis

Based on the computational efficiency analysis presented in Table 10, OnePiece exhibits superior hardware utilization and execution performance characteristics that support its viability for large-scale industrial deployment across both retrieval and ranking modes.

**Enhanced Hardware Utilization in Retrieval Mode.** OnePiece delivers significant efficiency gains over DLRM in retrieval tasks: a **25% reduction in inference time** (30ms vs. 40ms per request) while simultaneously achieving dramatic improvements in resource utilization. The **129% increase in Model FLOPs Utilization** (from 35% to 80%) indicates that OnePiece’s unified Transformer architecture exhibits superior compatibility with modern GPU parallelization paradigms, effectively leveraging tensor computation units that remain underutilized in traditional embedding-heavy architectures like DLRM. The concurrent 67% increase in memory utilization (from 30% to 50%) represents better resource exploitation, demonstrating more effective utilization of available hardware capacity while maintaining faster inference speeds. This efficiency improvement suggests that OnePiece’s architectural unification eliminates the computational bottlenecks inherent in DLRM’s heterogeneous component interactions, enabling more streamlined data flow and reduced memory transfer overhead.

Table 10 | Computational efficiency comparison between DLRM and OnePiece under both retrieval and ranking modes, evaluated on a single NVIDIA A30 GPU. Here, **MFU** denotes *Model FLOPs Utilization*, reflecting the ratio of achieved FLOPs to the theoretical peak, while **MU** denotes *Memory Utilization*, measured as the percentage of GPU memory capacity occupied during inference. The arrows indicate the preferred direction for each metric: “↑” for higher is better while “↓” for lower is better.

Retrieval Mode			
Method	Infer. Time↓	MFU↑	MU↑
DLRM	40ms/request	35%	30%
OnePiece	30ms/request (-25%)	80% (+129%)	50% (+67%)
Ranking Mode (batch size=128, KV-Cache enabled)			
Method	Infer. Time↓	MFU↑	MU↑
DLRM	109ms/batch	23%	29%
OnePiece (M=1)	110ms/batch (+0.9%)	67% (+191%)	38% (+31%)
OnePiece (M=4)	112ms/batch (+2.8%)	-	-
OnePiece (M=8)	115ms/batch (+5.5%)	-	-
OnePiece (M=12)	120ms/batch (+10.1%)	-	-

The superior hardware utilization translates directly to reduced operational costs, critical factors for large-scale industrial deployment where processing millions of requests daily.

**Controlled Computational Scaling in Ranking Mode.** The ranking mode evaluation reveals OnePiece’s capacity for efficient scaling with increased reasoning complexity. While OnePiece incurs modest overhead relative to DLRM (109ms baseline), the scaling behavior with block size exhibits desirable sub-linear characteristics: inference time increases from 110ms at  $M=1$  to 120ms at  $M=12$ , representing only a **10.1% overhead for 12× reasoning capacity expansion**. The progressive overhead pattern (0.9% → 2.8% → 5.5% → 10.1%) demonstrates that the block-wise reasoning mechanism achieves efficient computational amortization, where each additional reasoning block incurs diminishing marginal cost. This controlled scaling translates into a highly favorable efficiency-performance trade-off: as demonstrated in Table 6, C-AUC improves from 0.885 at  $M = 1$  to 0.927 at  $M = 12$ , representing a significant **4.7%** relative improvement, uncovering the potential for reasoning scaling in our unified framework. Particularly noteworthy is the dramatic **191% improvement in MFU** (from 23% to 67%) achieved even at  $M=1$ , indicating that OnePiece’s unified architecture aligns with GPU computational paradigms regardless of reasoning complexity. This efficiency gain stems from the infrastructure trick: the **KV-Caching** mechanism enables efficient batch processing while maintaining this superior hardware utilization, demonstrating the effectiveness in handling high-throughput production workloads. These findings establish OnePiece as a practical solution for production systems, offering configurable reasoning depth that enables flexible trade-offs between computational efficiency and model performance.

## 6. Related Work

**Context Engineering and Reasoning Enhancement in LLMs.** The Transformer architecture (Vaswani et al., 2017) has fundamentally transformed artificial intelligence across multiple domains (Brown et al., 2020; Dosovitskiy et al., 2020; Wang et al., 2017) and culminating in ChatGPT’s demonstration of emergent general-purpose capabilities in large language models (LLMs). Beyond established scaling laws (Chung et al., 2024; Isik et al., 2024; Zhao et al., 2023), current research has identified two orthogonal optimization pathways for maximizing LLM utility: **Context Engineering** and **Reason-**

**ing Enhancement.** In specific, context engineering encompasses the evolution from basic prompt engineering (Cain, 2024; Chen et al., 2023; Marvin et al., 2023) toward comprehensive context architecture (Amatriain, 2024; Mei et al., 2025; Velásquez-Henao et al., 2023) that orchestrate prompts, external knowledge bases, persistent memory, and interaction histories into input context to effectively steer reasoning and generation. The reasoning enhancement pathway advances computational reasoning through sophisticated methodologies including chain-of-thought prompting (Diao et al., 2023; Wei et al., 2022; Zhang et al., 2022), structured tree-based (Cao et al., 2023; Long, 2023; Yao et al., 2023a) and graph-based reasoning (Besta et al., 2024; Jin et al., 2024; Yao et al., 2023b), and iterative self-refinement mechanisms (Madaan et al., 2023; Ranaldi and Freitas, 2024; Yan et al., 2024) that augment latent multi-step reasoning processes (Hao et al., 2024; Su et al., 2025; Zhu et al., 2025), enabling logical decomposition and strategic planning that transcends statistical pattern matching. The convergence of these complementary research directions transforms LLMs from probabilistic sequence generators into adaptive reasoning architectures capable of complex problem-solving and knowledge integration.

**Context Engineering and Reasoning Enhancement in Ranking Systems.** The evolution of ranking systems has progressively integrated advances in language modeling, with early explorations mainly focusing on architectural adaptations by directly applying Transformer-based models like SASRec (Kang and McAuley, 2018) and BERT4Rec (Sun et al., 2019). Inspired by advances in large language models (LLMs), recent research has crystallized into two promising research directions that mirror established LLM optimization pathways. First, context engineering-based approaches focus on developing sophisticated user sequence representations through dynamic behavior modeling and retrieval-augmented methods, with representative works (Li et al., 2024; Lin et al., 2024) transitioning from conventional recency-based sequence modeling toward semantic behavior retrieval algorithms that extract contextually relevant historical interactions, demonstrating enhanced robustness and improved cold-start prediction accuracy. Second, reasoning enhancement-based methods have bifurcated into explicit and implicit paradigms (Wang et al., 2024; Yu et al., 2025; Zhang et al., 2025a), where explicit reasoning approaches (Fang et al., 2025; Gu et al., 2025) employ multi-expert pipelines and reflection-refinement mechanisms to generate explainable intermediate reasoning at the cost of computational overhead, while implicit reasoning methods (Liu et al., 2025; Tang et al., 2025; Zhang et al., 2025b) leverage continuous hidden state autoregression and depth-recurrent architectures to achieve enhanced reasoning capabilities without explicit text generation requirements, offering computational efficiency while maintaining reasoning depth through specialized alignment objectives.

While existing efforts have explored various LLM-inspired techniques for ranking systems, they lack a principled approach that captures the fundamental design principles driving LLM success. To this end, we propose OnePiece, which addresses this gap by establishing a unified architectural paradigm that comprehensively adapts the core design philosophy of LLMs to ranking systems.

## 7. Conclusion and Future Directions

In this paper, we introduce **OnePiece**, a unified framework that brings structured context engineering and block-wise latent reasoning into industrial cascaded ranking systems. Built upon a pure Transformer backbone, OnePiece incorporates context engineering to organize heterogeneous signals (interaction history, preference anchors, situational descriptors, and candidate items) into a structured tokenized sequence, equips the model with multi-step reasoning capacity, and optimizes this process through a progressive multi-task training strategy that leverages naturally available user feedback chains. Extensive offline experiments validate the effectiveness of each design, reveal favorable scaling with respect to preference anchor length, training data span, and block size. Online A/B testing at Shopee’s main personalized search scenario demonstrates that OnePiece achieves consistent gains

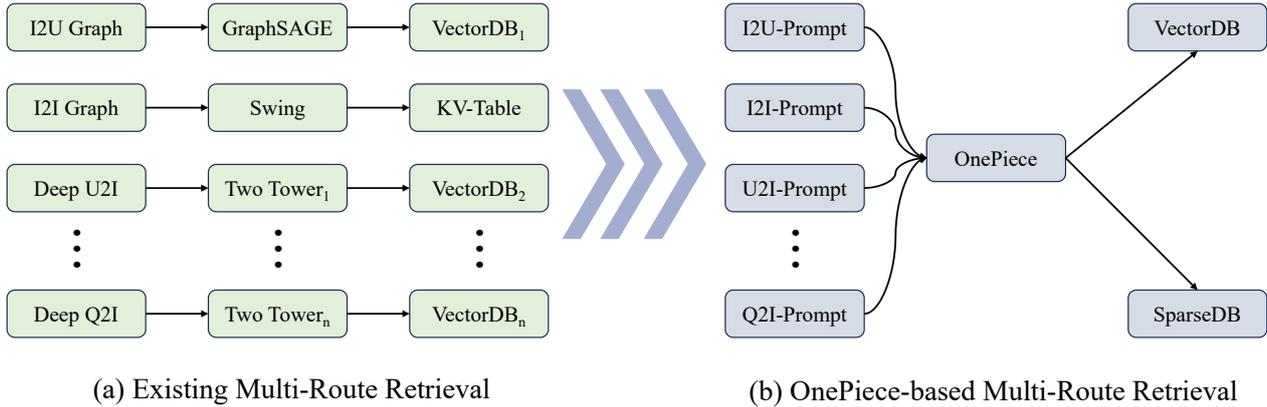


Figure 7 | Comparison between existing multi-route retrieval systems and OnePiece-based unified architecture. (a) Traditional multi-route retrieval requires maintaining separate models with distinct parameters for different retrieval pathways (I2U, I2I, U2I, Q2I, etc.), each utilizing specialized architectures and storage systems. (b) OnePiece achieves unified multi-route retrieval through a single model that processes tailored prompts for different retrieval scenarios, enabling “One For All” functionality while reducing system complexity and maintenance overhead.

across diverse business metrics, including over +2% GMV/UU improvement and +2.90% advertising revenue improvement, while also exhibiting superior efficiency and hardware utilization. These results position OnePiece as a promising new paradigm for building scalable, reasoning-driven ranking models in real-world industrial environments.

Looking ahead, several promising research directions emerge from this work:

- **Unified Multi-Route Retrieval:** As illustrated in Fig. 7, existing multi-route retrieval systems require distinct design principles, training data, and model parameters for each retrieval pathway to achieve diversified personalization and comprehensive coverage of user interests. However, maintaining multiple heterogeneous models is resource-intensive and prone to redundancy across different routes. Our unified OnePiece architecture paves the way for “**One For All**” multi-route retrieval, where a single unified model can serve diverse recommendation objectives through tailored context engineering (e.g., adapting interaction history, preference anchors, and situational descriptors based on specific scenario characteristics). The empirical evidence from Sec. 5.3 demonstrates strong diversity and exclusivity performance, supporting the feasibility of streamlined multi-route system design in industrial applications.
- **Scalable Latent Reasoning:** While this work represents the first successful deployment of latent reasoning in industrial-scale ranking systems through multi-task progressive supervision, this approach also reveals inherent limitations in reasoning scalability. The challenge lies in obtaining sufficient multi-task signals to effectively supervise intermediate reasoning processes, which constrains our ability to scale reasoning capabilities further. Future research should explore more effective scaling methodologies for latent reasoning, such as incorporating online user feedback through reinforcement learning to adaptively determine optimal reasoning depth, or developing organic integration between model self-exploration and multi-task supervision processes that enable autonomous reasoning evolution.

## References

- Q. Ai, D. N. Hill, S. Vishwanathan, and W. B. Croft. A zero attention model for personalized product search. In Proceedings of the 28th ACM International Conference on Information and Knowledge Management, pages 379–388, 2019.
- X. Amatriain. Prompt design and engineering: Introduction and advanced methods. arXiv preprint arXiv:2401.14423, 2024.
- Y. Bengio, J. Louradour, R. Collobert, and J. Weston. Curriculum learning. In Proceedings of the 26th annual international conference on machine learning, pages 41–48, 2009.
- M. Besta, N. Blach, A. Kubicek, R. Gerstenberger, M. Podstawski, L. Gianinazzi, J. Gajda, T. Lehmann, H. Niewiadomski, P. Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In Proceedings of the AAAI conference on artificial intelligence, volume 38, pages 17682–17690, 2024.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- W. Cain. Prompting change: Exploring prompt engineering in large language model ai and its potential to transform education. TechTrends, 68(1):47–57, 2024.
- S. Cao, J. Zhang, J. Shi, X. Lv, Z. Yao, Q. Tian, J. Li, and L. Hou. Probabilistic tree-of-thought reasoning for answering knowledge-intensive complex questions. arXiv preprint arXiv:2311.13982, 2023.
- B. Chen, Z. Zhang, N. Langrené, and S. Zhu. Unleashing the potential of prompt engineering in large language models: a comprehensive review. arXiv preprint arXiv:2310.14735, 2023.
- B. Chen, X. Guo, S. Wang, Z. Liang, Y. Lv, Y. Ma, X. Xiao, B. Xue, X. Zhang, Y. Yang, et al. Onesearch: A preliminary exploration of the unified end-to-end generative framework for e-commerce search. arXiv preprint arXiv:2509.03236, 2025.
- H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, Y. Li, X. Wang, M. Dehghani, S. Brahma, et al. Scaling instruction-finetuned language models. Journal of Machine Learning Research, 25(70):1–53, 2024.
- J. Deng, S. Wang, K. Cai, L. Ren, Q. Hu, W. Ding, Q. Luo, and G. Zhou. Onerec: Unifying retrieve and rank with generative recommender and iterative preference alignment. arXiv preprint arXiv:2502.18965, 2025.
- S. Diao, P. Wang, Y. Lin, R. Pan, X. Liu, and T. Zhang. Active prompting with chain-of-thought for large language models. arXiv preprint arXiv:2302.12246, 2023.
- Q. Dong, L. Li, D. Dai, C. Zheng, J. Ma, R. Li, H. Xia, J. Xu, Z. Wu, B. Chang, et al. A survey on in-context learning. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 1107–1128, 2024.
- A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.

- Y. Fang, W. Wang, Y. Zhang, F. Zhu, Q. Wang, F. Feng, and X. He. Reason4rec: Large language models for recommendation with deliberative user preference alignment. [arXiv preprint arXiv:2502.02061](#), 2025.
- C. Fu, K. Wang, J. Wu, Y. Chen, G. Huzhang, Y. Ni, A. Zeng, and Z. Zhou. Residual multi-task learner for applied ranking. In [Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining](#), pages 4974–4985, 2024.
- H. Gu, R. Zhong, Y. Xia, W. Yang, C. Lu, P. Jiang, and K. Gai. R4ec: A reasoning, reflection, and refinement framework for recommendation systems. [arXiv preprint arXiv:2507.17249](#), 2025.
- S. Hao, S. Sukhbaatar, D. Su, X. Li, Z. Hu, J. Weston, and Y. Tian. Training large language models to reason in a continuous latent space. [arXiv preprint arXiv:2412.06769](#), 2024.
- J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. *iee*. In [CVF Conference on Computer Vision and Pattern Recognition](#). IEEE, pages 7132–41, 2017.
- J. Huang and K. C.-C. Chang. Towards reasoning in large language models: A survey. In [Findings of the Association for Computational Linguistics: ACL 2023](#), pages 1049–1065, 2023.
- P.-S. Huang, X. He, J. Gao, L. Deng, A. Acero, and L. Heck. Learning deep structured semantic models for web search using clickthrough data. In [Proceedings of the 22nd ACM international conference on Information & Knowledge Management](#), pages 2333–2338, 2013.
- P. Indyk and R. Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In [Proceedings of the thirtieth annual ACM symposium on Theory of computing](#), pages 604–613, 1998.
- B. Isik, N. Ponomareva, H. Hazimeh, D. Paparas, S. Vassilvitskii, and S. Koyejo. Scaling laws for downstream task performance of large language models. In [ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models](#), 2024.
- B. Jin, C. Xie, J. Zhang, K. K. Roy, Y. Zhang, Z. Li, R. Li, X. Tang, S. Wang, Y. Meng, et al. Graph chain-of-thought: Augmenting large language models by reasoning on graphs. [arXiv preprint arXiv:2404.07103](#), 2024.
- W.-C. Kang and J. McAuley. Self-attentive sequential recommendation. In [2018 IEEE international conference on data mining \(ICDM\)](#), pages 197–206. IEEE, 2018.
- Y. Li, J. Wang, T. Dai, J. Zhu, J. Yuan, R. Zhang, and S.-T. Xia. Rat: Retrieval-augmented transformer for click-through rate prediction. In [Companion Proceedings of the ACM Web Conference 2024](#), pages 867–870, 2024.
- J. Lin, R. Shan, C. Zhu, K. Du, B. Chen, S. Quan, R. Tang, Y. Yu, and W. Zhang. Rella: Retrieval-enhanced large language models for lifelong sequential behavior comprehension in recommendation. In [Proceedings of the ACM Web Conference 2024](#), pages 3497–3508, 2024.
- J. Lin, X. Dai, Y. Xi, W. Liu, B. Chen, H. Zhang, Y. Liu, C. Wu, X. Li, C. Zhu, et al. How can recommender systems benefit from large language models: A survey. [ACM Transactions on Information Systems](#), 43(2):1–47, 2025.
- E. Liu, B. Zheng, X. Wang, W. X. Zhao, J. Wang, S. Chen, and J.-R. Wen. Lares: Latent reasoning for sequential recommendation. [arXiv preprint arXiv:2505.16865](#), 2025.
- J. Long. Large language model guided tree-of-thought. [arXiv preprint arXiv:2305.08291](#), 2023.

- A. Madaan, N. Tandon, P. Gupta, S. Hallinan, L. Gao, S. Wiegrefe, U. Alon, N. Dziri, S. Prabhunoye, Y. Yang, et al. Self-refine: Iterative refinement with self-feedback. Advances in Neural Information Processing Systems, 36:46534–46594, 2023.
- Y. A. Malkov and D. A. Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. IEEE transactions on pattern analysis and machine intelligence, 42(4):824–836, 2018.
- G. Marvin, N. Hellen, D. Jjingo, and J. Nakatumba-Nabende. Prompt engineering in large language models. In International conference on data intelligence and cognitive informatics, pages 387–402. Springer, 2023.
- L. Mei, J. Yao, Y. Ge, Y. Wang, B. Bi, Y. Cai, J. Liu, M. Li, Z.-Z. Li, D. Zhang, et al. A survey of context engineering for large language models. arXiv preprint arXiv:2507.13334, 2025.
- M. Naumov, D. Mudigere, H.-J. M. Shi, J. Huang, N. Sundaraman, J. Park, X. Wang, U. Gupta, C.-J. Wu, A. G. Azzolini, et al. Deep learning recommendation model for personalization and recommendation systems. arXiv preprint arXiv:1906.00091, 2019.
- A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PmLR, 2021.
- L. Ranaldi and A. Freitas. Self-refine instruction-tuning for aligning reasoning in language models. arXiv preprint arXiv:2405.00402, 2024.
- Z. Shen, H. Yan, L. Zhang, Z. Hu, Y. Du, and Y. He. Codi: Compressing chain-of-thought into continuous space via self-distillation. arXiv preprint arXiv:2502.21074, 2025.
- D. Su, H. Zhu, Y. Xu, J. Jiao, Y. Tian, and Q. Zheng. Token assorted: Mixing latent and text tokens for improved language model reasoning. arXiv preprint arXiv:2502.03275, 2025.
- F. Sun, J. Liu, J. Wu, C. Pei, X. Lin, W. Ou, and P. Jiang. Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer. In Proceedings of the 28th ACM international conference on information and knowledge management, pages 1441–1450, 2019.
- J. Tang, S. Dai, T. Shi, J. Xu, X. Chen, W. Chen, J. Wu, and Y. Jiang. Think before recommend: Unleashing the latent reasoning power for sequential recommendation. arXiv preprint arXiv:2503.22675, 2025.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
- J. D. Velásquez-Henao, C. J. Franco-Cardona, and L. Cadavid-Higueta. Prompt engineering: a methodology for optimizing interactions with ai-language models in the field of engineering. Dyna, 90(SPE230):9–17, 2023.
- R. Wang, R. Shivanna, D. Cheng, S. Jain, D. Lin, L. Hong, and E. Chi. Dcn v2: Improved deep & cross network and practical lessons for web-scale learning to rank systems. In Proceedings of the web conference 2021, pages 1785–1797, 2021.
- Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, et al. Tacotron: Towards end-to-end speech synthesis. arXiv preprint arXiv:1703.10135, 2017.

- Y. Wang, C. Tian, B. Hu, Y. Yu, Z. Liu, Z. Zhang, J. Zhou, L. Pang, and X. Wang. Can small language models be good reasoners for sequential recommendation? In Proceedings of the ACM Web Conference 2024, pages 3876–3887, 2024.
- J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.
- L. Wu, Z. Zheng, Z. Qiu, H. Wang, H. Gu, T. Shen, C. Qin, C. Zhu, H. Zhu, Q. Liu, et al. A survey on large language models for recommendation. World Wide Web, 27(5):60, 2024.
- H. Yan, Q. Zhu, X. Wang, L. Gui, and Y. He. Mirror: A multiple-perspective self-reflection method for knowledge-rich reasoning. arXiv preprint arXiv:2402.14963, 2024.
- S. Yao, D. Yu, J. Zhao, I. Shafran, T. Griffiths, Y. Cao, and K. Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. Advances in neural information processing systems, 36:11809–11822, 2023a.
- Y. Yao, Z. Li, and H. Zhao. Beyond chain-of-thought, effective graph-of-thought reasoning in language models. arXiv preprint arXiv:2305.16582, 2023b.
- Q. Yu, K. Fu, S. Zhang, Z. Lv, F. Wu, and F. Wu. Thinkrec: Thinking-based recommendation via llm. arXiv preprint arXiv:2505.15091, 2025.
- J. Zhai, L. Liao, X. Liu, Y. Wang, R. Li, X. Cao, L. Gao, Z. Gong, F. Gu, J. He, et al. Actions speak louder than words: trillion-parameter sequential transducers for generative recommendations. In Proceedings of the 41st International Conference on Machine Learning, pages 58484–58509, 2024.
- J. Zhang, B. Zhang, W. Sun, H. Lu, W. X. Zhao, Y. Chen, and J.-R. Wen. Slow thinking for sequential recommendation. arXiv preprint arXiv:2504.09627, 2025a.
- Y. Zhang, W. Xu, X. Zhao, W. Wang, F. Feng, X. He, and T.-S. Chua. Reinforced latent reasoning for llm-based recommendation. arXiv preprint arXiv:2505.19092, 2025b.
- Z. Zhang, A. Zhang, M. Li, and A. Smola. Automatic chain of thought prompting in large language models. arXiv preprint arXiv:2210.03493, 2022.
- W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, et al. A survey of large language models. arXiv preprint arXiv:2303.18223, 1(2), 2023.
- G. Zhou, X. Zhu, C. Song, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, and K. Gai. Deep interest network for click-through rate prediction. In Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining, pages 1059–1068, 2018.
- R.-J. Zhu, T. Peng, T. Cheng, X. Qu, J. Huang, D. Zhu, H. Wang, K. Xue, X. Zhang, Y. Shan, et al. A survey on latent reasoning. arXiv preprint arXiv:2507.06203, 2025.
- Y. Zhu, H. Yuan, S. Wang, J. Liu, W. Liu, C. Deng, H. Chen, Z. Liu, Z. Dou, and J.-R. Wen. Large language models for information retrieval: A survey. arXiv preprint arXiv:2308.07107, 2023.

## Appendix

### A. Offline Dataset Construction

**Retrieval Stage.** The objective of the retrieval task is to retrieve items that users may potentially interact with in the future. Therefore, we focus on learning impression and click objectives. For each sample, we first filter out session request data where users did not exhibit click behaviors from the online data. Then, we construct the training samples as follows: (1) the  $m$  clicked items serve as positive samples for both impression and click tasks; (2) the  $n$  exposed but unclicked items serve as positive samples for the impression task while simultaneously acting as negative samples for the click task; (3) we further sample  $k$  items from unexposed items within the top-500 results from the ranking stage as additional negative samples; (4) we sample  $l$  items from the same category as the clicked items, which serve as hard negative samples to enhance model convergence and mitigate homogeneous recommendation risks. The specific values of  $m$ ,  $n$ ,  $k$ , and  $l$  are determined based on domain expert experience and empirical validation.

**Ranking Stage.** As the downstream stage following retrieval, the ranking task requires further refinement of user interests and preferences to more accurately compute relevance scores for candidate items. We filter session request data to retain only those with click behaviors. For training sample construction, we use task-specific interaction types (*i.e.*, impression, click, add-to-cart, order) as positive samples, while interactions from preceding tasks in the conversion funnel serve as negative samples for each respective task (*e.g.*, for the order prediction task, exposed, clicked, and add-to-cart items that were not purchased serve as negative samples). Additionally, analogous to the retrieval stage, we enhance the negative sample pool by randomly sampling items from the top-500 ranking results that were not exposed to users, which serve as augmented hard negative samples to improve model performance.

### B. Context Engineering Details

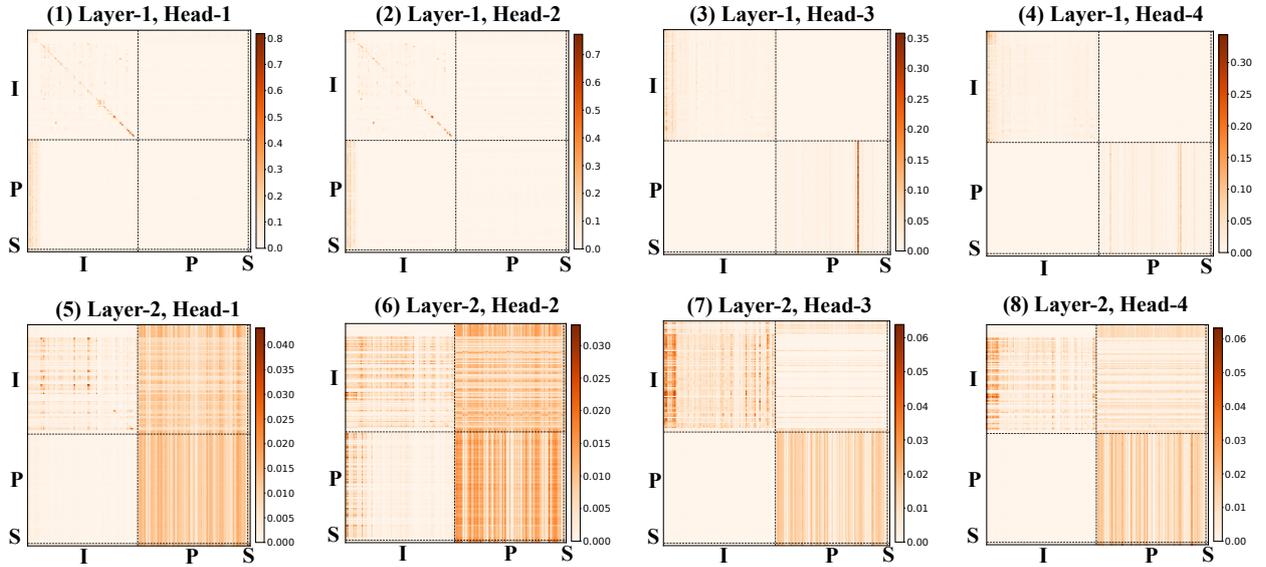
For the input components, namely Interaction History (IH), Preference Anchor (PA), Situational Descriptors (SD), and Candidate Item Set (CIS), we detail the construction as follows. As CIS involves straightforward item features, here, we mainly focus on the other three components.

**Interaction History (IH):** We utilize users' click sequences, shopping cart sequences, and purchase sequences from the most recent month, chronologically sorted and merged. Each item in the sequence encompasses features including item ID, category, shop, and other relevant attributes.

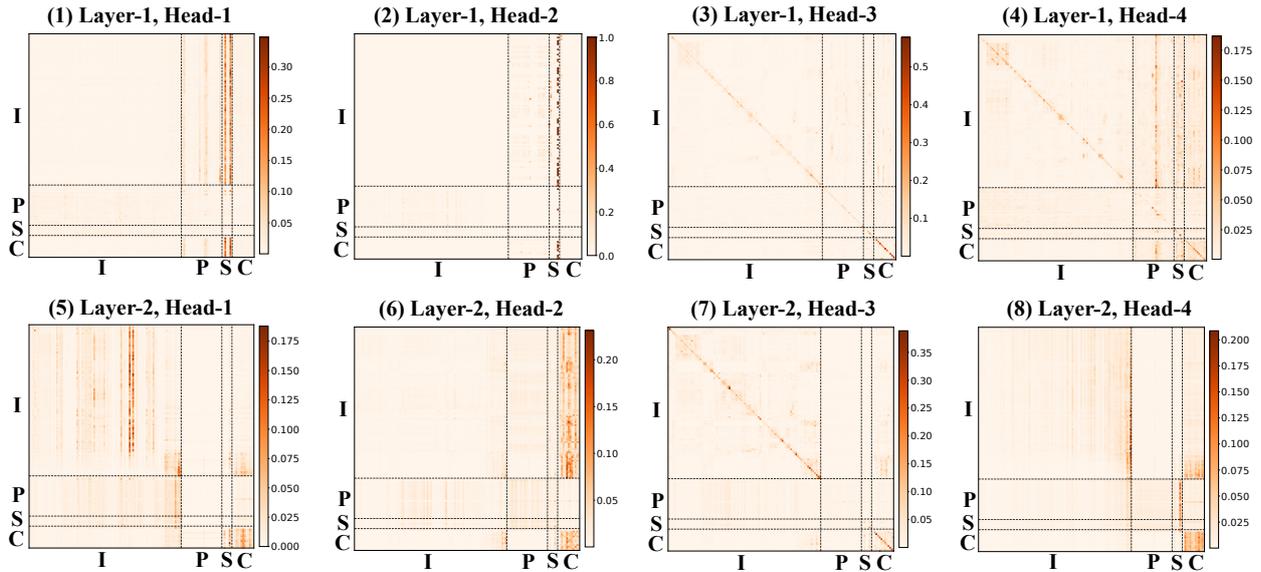
**Preference Anchor (PA):** We aggregate users' query-associated top-k exposures, top-k clicks, and top-k purchases, utilizing [BOS] and [EOS] tokens to separate different sequence types. Unlike the IH component, we employ sequence concatenation rather than mixing approaches.

**Situational Descriptors (SD):** This component covers user profiles and query-related contextual information. We introduce two specialized tokens: User Token (integrating user ID, age, location, and other demographic features) and Query Token (incorporating query ID, textual content, query popularity, and related information).

### C. Attention Visualization Analysis



(a) Case Studies on OnePiece Attention Analysis (Retrieval Mode).



(b) Case Studies on OnePiece Attention Analysis (Ranking Mode).

Figure 8 | OnePiece Attention Analysis in Different Modes. The attention maps visualize the attention weights between different input components: Interaction History (I), Preference Anchor (P), Situational Descriptor (S), and Candidate Item Set (C). In these visualizations, the **y-axis** represents the Query, while the **x-axis** represents the Key-Value, corresponding to the attention weight matrix commonly used in Transformer-like architectures.

**Attention Analysis of Context Input.** Figure 8 illustrates how our unified token sequence (IH, PA, SD, CIS) facilitates structured preference reasoning across both retrieval and ranking stages. Across modes, we observe two consistent trends. First, *layer-wise evolution* is evident: Layer-1 heads

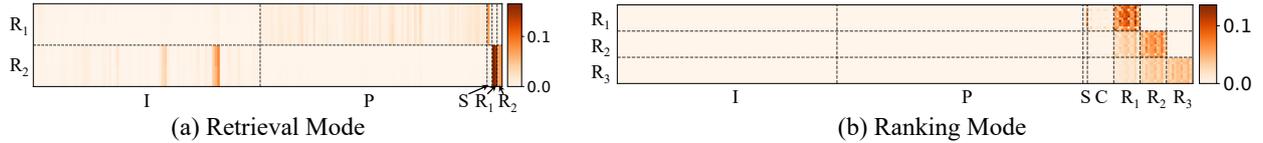


Figure 9 | Attention visualization of multi-step block-wise reasoning in OnePiece. The heatmaps show attention weights between reasoning blocks (y-axis, as queries) and input components with previous reasoning outputs (x-axis, as keys and values) for **(a)** retrieval mode with two reasoning steps ( $R_1$ ,  $R_2$ ) and **(b)** ranking mode with three reasoning steps ( $R_1$ ,  $R_2$ ,  $R_3$ ). Context Tokens include Interaction History (I), Preference Anchors (P), Situational Descriptors (S), and Candidate Items (C, ranking only).

(e.g., Figure 8(a)-1–4 and Figure 8(b)-1–3) primarily form concentrated or diagonal attention bands, highlighting localized sequential processing within IH tokens or short-span links between SD and CIS. In contrast, Layer-2 heads (Figure 8(a)-5–8; Figure 8(b)-5–8) develop multi-region attention patterns that simultaneously connect multiple token groups, suggesting a transition from localized to global integration. Second, *head-level specialization* emerges clearly. Within the same layer, different heads adopt complementary strategies: some emphasize intra-component coherence (e.g., IH–IH diagonals in Figure 8(b)-7), while others allocate strong weights to cross-component flows (e.g., SD→IH in Figure 8(a)-5,6; CIS→IH in Figure 8(b)-6). Together, these patterns validate that the model learns hierarchical and diversified reasoning strategies rather than redundantly replicating attention across different heads of OnePiece.

Beyond these shared behaviors, each mode exhibits distinctive characteristics aligned with its task. In retrieval mode (Figure 8(a)), the three-token design (IH, PA, SD) fosters *structured and compact cross-component attention*. For example, Head-1 and Head-2 at Layer-2 (Figure 8(a)-5,6) concentrate on linking IH with PA, highlighting how anchors guide long-term preference recall, while Head-4 (Figure 8(a)-8) reinforces SD→IH connections, grounding retrieval in situational context. These interactions remain relatively localized, reflecting retrieval’s objective of coarse-grained candidate filtering. By contrast, ranking mode (Figure 8(b)) introduces CIS tokens, fundamentally expanding the attention space to *four-way interactions*. This is most pronounced in Head-2 and Head-4 at Layer-2 (Figure 8(b)-6,8), where attention flows span IH, PA, SD, and CIS simultaneously, enabling joint evaluation of user preference signals against explicit candidate items. Importantly, we see *functional differentiation*: IH tokens preserve temporal sequentiality (Figure 8(b)-7), while CIS tokens actively integrate with PA and SD to enable fine-grained candidate comparison (Figure 8(b)-4). These entangled attention maps underscore the ranking stage’s role in nuanced discrimination among candidates, complementing the retrieval stage’s more selective filtering focus.

**Attention Analysis of Multi-Step Block-wise Reasoning.** The attention patterns in Figure 9 demonstrate how reasoning blocks progressively query different information sources to refine user representations. In **retrieval mode**,  $R_1$  exhibits strong concentrated attention on situational descriptors (S) and moderate attention on preference anchors (P), with minimal attention to interaction history (I). This pattern indicates that initial reasoning prioritizes contextual and query-specific signals for user intent understanding.  $R_2$  shows a pivotal shift, developing concentrated attention on specific regions within interaction history (I) while also incorporating information from the previous reasoning block  $R_1$ . This evolution from situational-preference focus to selective behavioral pattern recognition demonstrates how progressive reasoning enables the model to transition from broad contextual understanding to targeted sequential preference extraction.

In **ranking mode**, the three-step reasoning process ( $R_1 \rightarrow R_2 \rightarrow R_3$ ) reveals increasingly sophisticated attention integration with a notable hierarchical information enhancement pattern. As reasoning progresses, later blocks demonstrate stronger attention to more recent reasoning outputs, with  $R_3$  showing pronounced attention to  $R_2$  while exhibiting relatively weaker attention to  $R_1$ . This attention pattern suggests that each reasoning step progressively consolidates and refines information from previous steps, with more recent reasoning blocks containing higher-level abstractions that subsume earlier insights. The model effectively learns to prioritize the most refined representations rather than repeatedly accessing raw earlier outputs, indicating an efficient information compression mechanism where each reasoning step builds upon increasingly compressed preference understanding to achieve discriminative candidate evaluation.